

# Heterogeneous Data-Centric Architectures for Modern Data-Intensive Applications: Case Studies in Machine Learning and Databases

P&S Processing-in-Memory  
Fall 2022  
17 January 2023

**Geraldo F. Oliveira**

**SAFARI**

**ETH** zürich

# Google Neural Network Models for Edge Devices: Analyzing and Mitigating Machine Learning Inference Bottlenecks

P&S Processing-in-Memory  
Fall 2022  
17 January 2023

**Amirali Boroumand**

**Saugata Ghose**

**Berkin Akin**

**Ravi Narayanaswami**

**Geraldo F. Oliveira**

**Xiaoyu Ma**

**Eric Shiu**

**Onur Mutlu**

**SAFARI**

# Executive Summary

**Context:** We extensively analyze a state-of-the-art edge ML accelerator (Google Edge TPU) using 24 Google edge models

- Wide range of models (CNNs, LSTMs, Transducers, RCNNs)

**Problem:** The Edge TPU accelerator suffers from **three challenges:**

- It operates **significantly below** its peak throughput
- It operates **significantly below** its theoretical energy efficiency
- It **inefficiently** handles memory accesses

**Key Insight:** These shortcomings arise from **the monolithic design** of the Edge TPU accelerator

- The Edge TPU accelerator design does not account for **layer heterogeneity**

**Key Mechanism:** A new framework called **Mensa**

- Mensa consists of heterogeneous accelerators whose dataflow and hardware are specialized for specific families of layers

**Key Results:** We design a version of Mensa for Google edge ML models

- Mensa improves performance and energy by **3.0X** and **3.1X**
- Mensa reduces cost and improves area efficiency

# Outline

**1** Introduction

**2** Edge TPU and Model Characterization

**3** Mensa Framework

**4** Mensa-G: Mensa for Google Edge Models

**5** Evaluation

**6** Conclusion

# Outline

## 1 Introduction

## 2 Edge TPU and Model Characterization

## 3 Mensa Framework

## 4 Mensa-G: Mensa for Google Edge Models

## 5 Evaluation

## 6 Conclusion

# Why ML on Edge Devices?

Significant interest in pushing ML inference computation directly to edge devices



**Privacy**



**Connectivity**



**Latency**



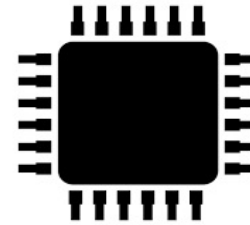
**Bandwidth**

# Why Specialized ML Accelerator?

Edge devices have limited battery and computation budget

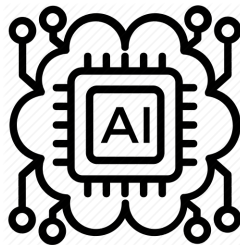


**Limited** Power Budget



**Limited** Computational Resources

Specialized accelerators can significantly improve inference latency and energy consumption

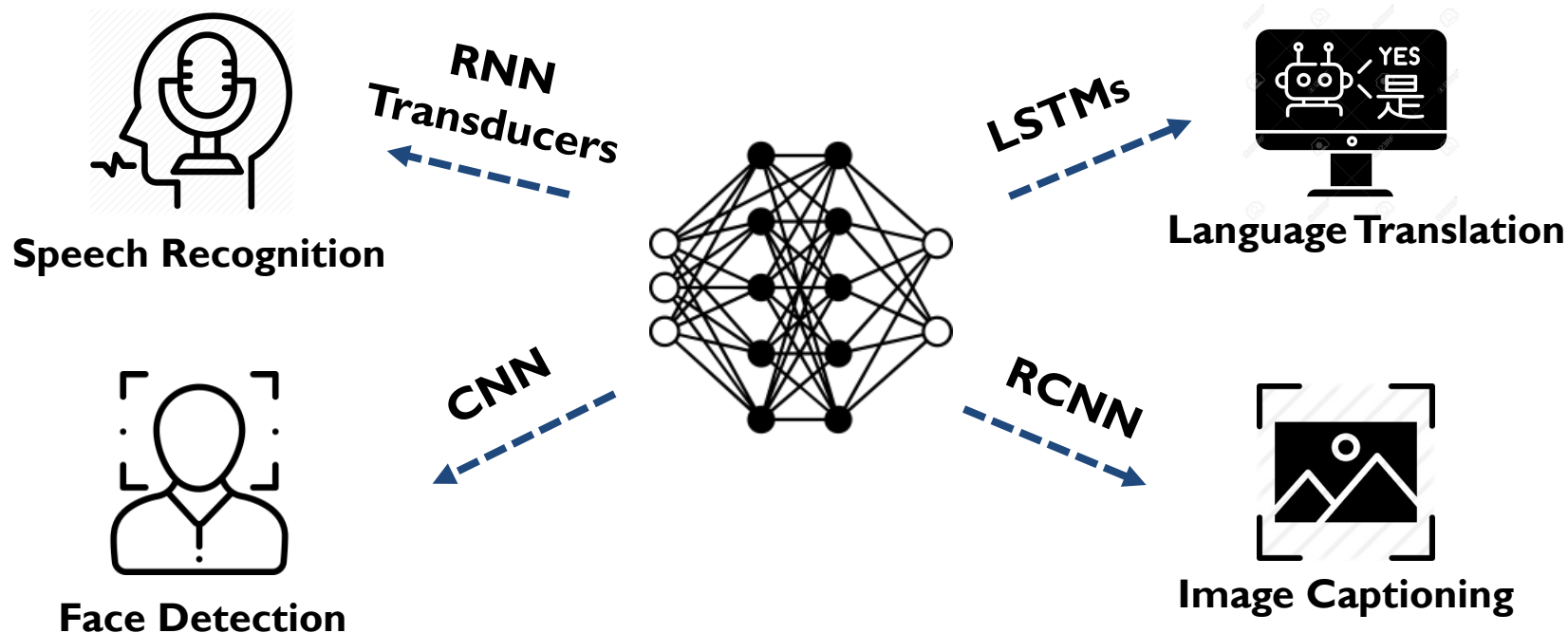


**Apple Neural Engine (A12)**



**Google Edge TPU**

# Myriad of Edge Neural Network Models



**Challenge: edge ML accelerators have to execute inference efficiently across a wide variety of NN models**



# Outline

1 Introduction

**2 Edge TPU and Model Characterization**

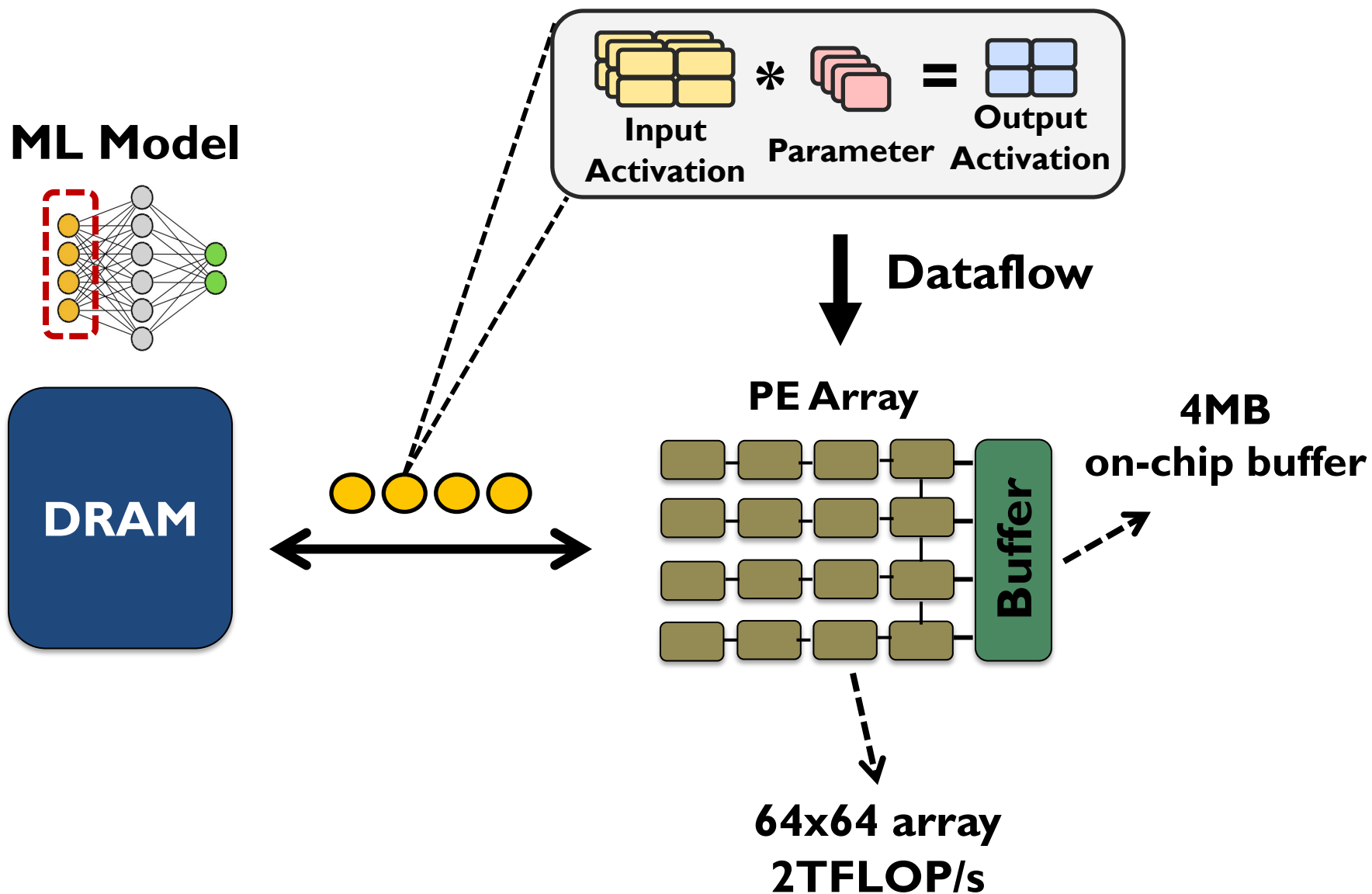
3 Mensa Framework

4 Mensa-G: Mensa for Google Edge Models

5 Evaluation

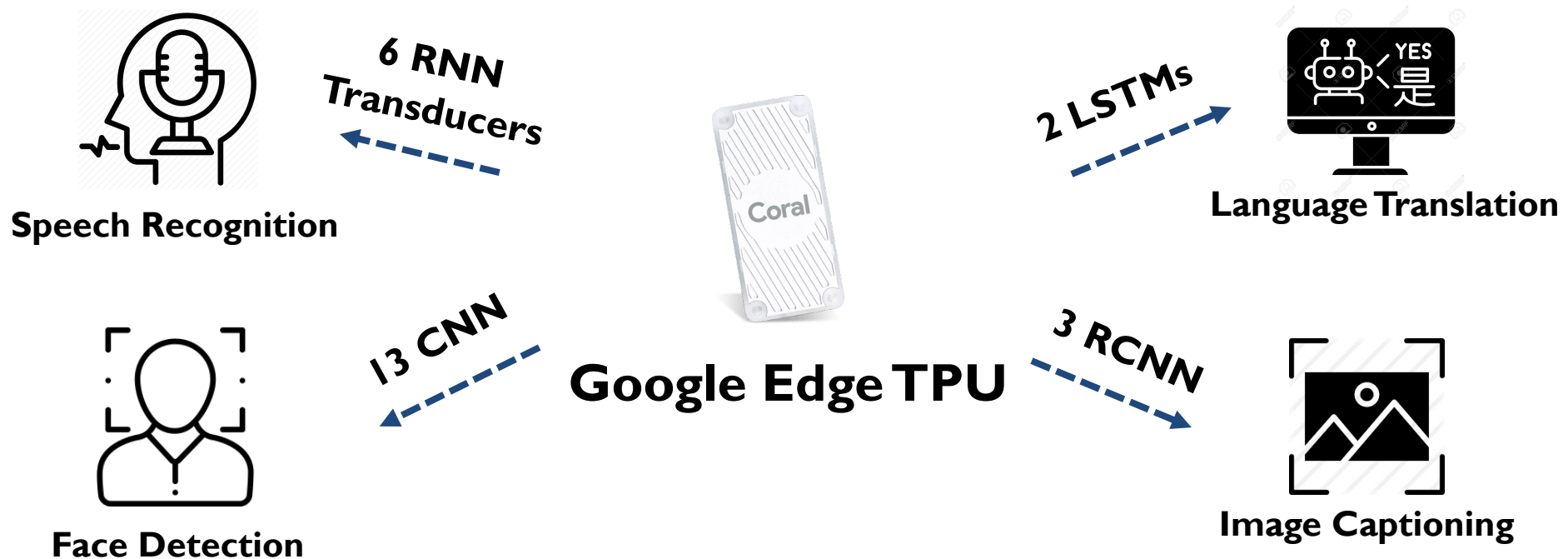
6 Conclusion

# Edge TPU: Baseline Accelerator



# Google Edge NN Models

## We analyze inference execution using 24 edge NN models



# Major Edge TPU Challenges

We find that the accelerator suffers from three major challenges:

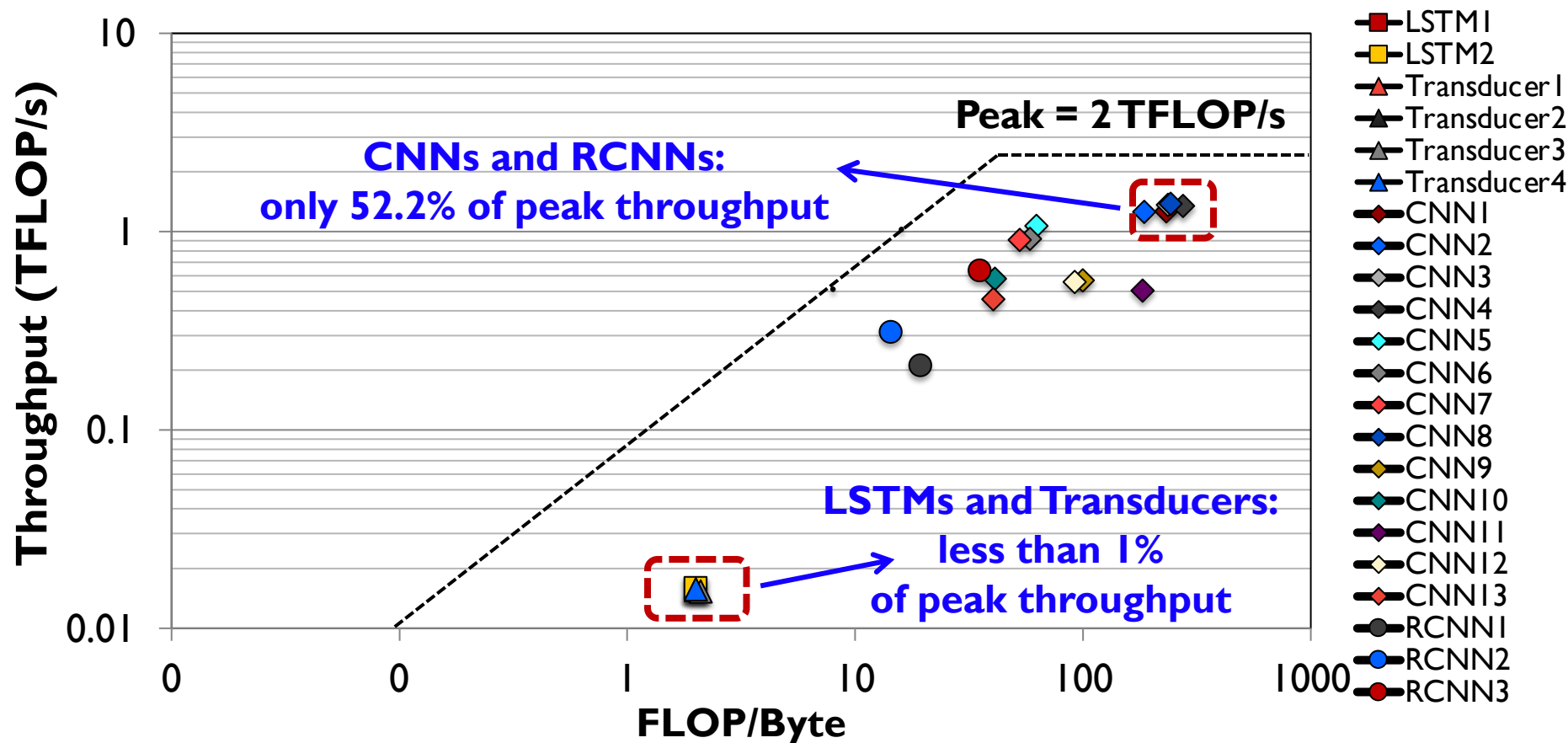
1 Operates significantly below its peak throughput

2 Operates significantly below its peak energy efficiency

3 Handles memory accesses inefficiently

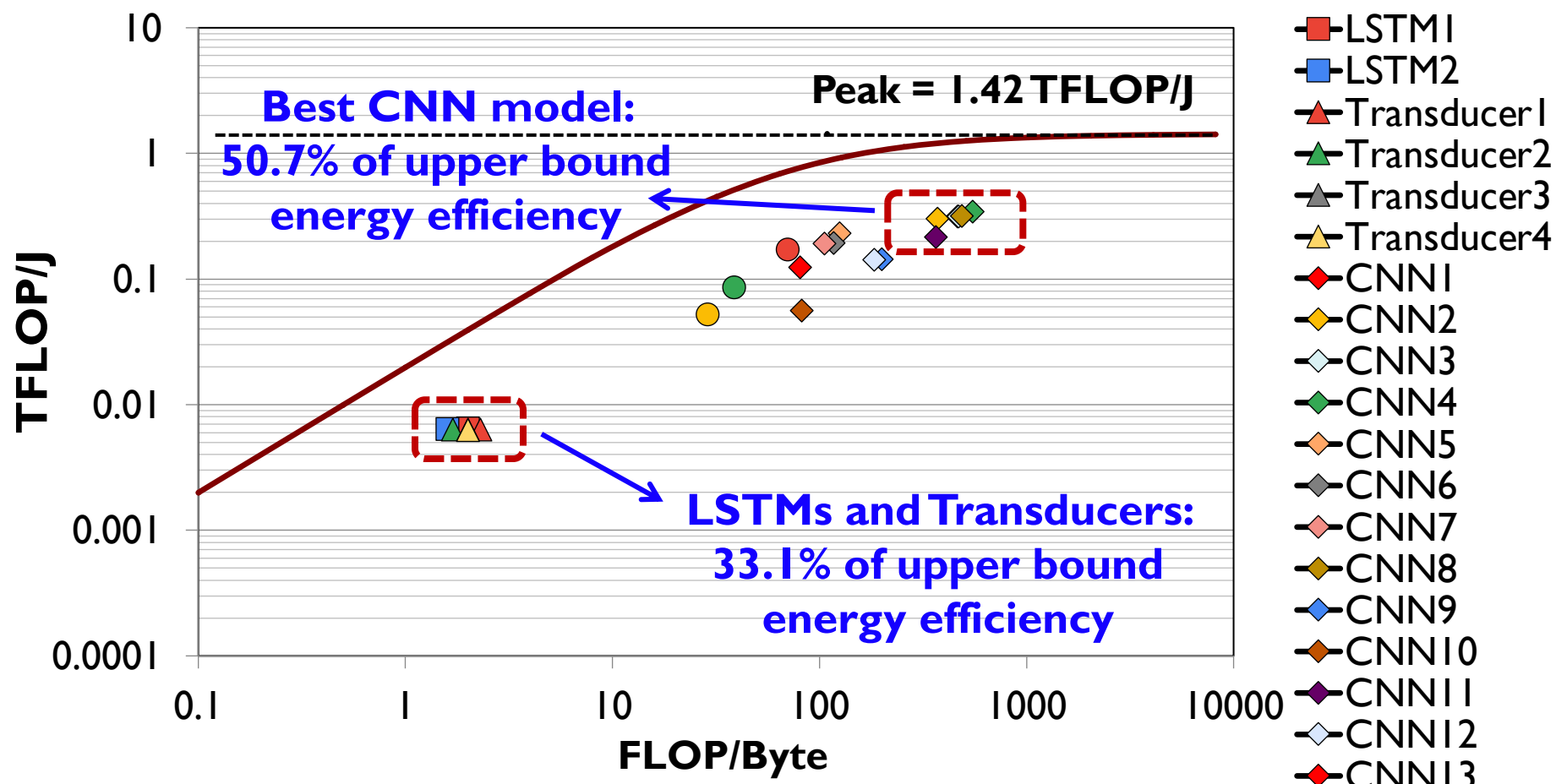
# (I) High Resource Underutilization

We find that the accelerator operates significantly below its peak throughput across all models



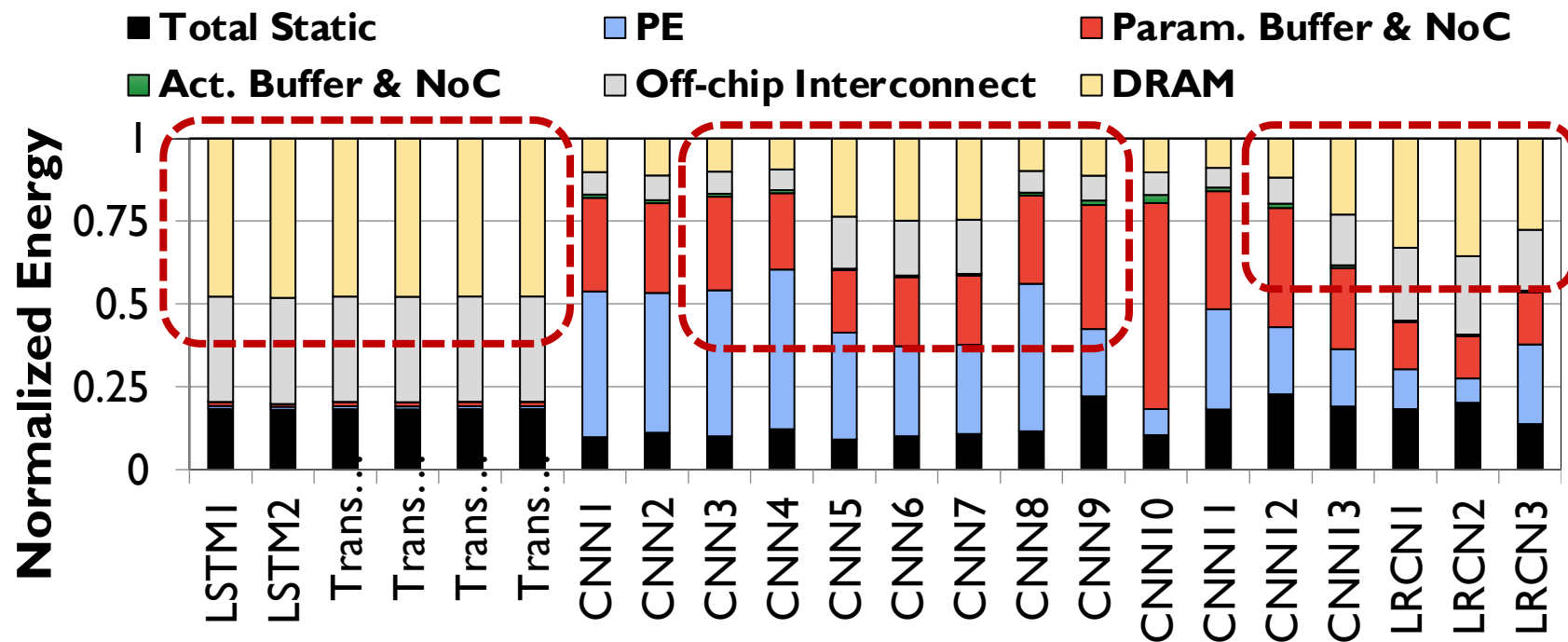
## (2) Low Energy Efficiency

The accelerator operates far below its upper bound energy efficiency



# (3) Inefficient Memory Access Handling

Parameter traffic (off-chip and on-chip) takes a large portion of the inference energy and performance



**46%** and **31%** of total energy goes to **off-chip parameter traffic** and **distributing parameters** across PE array

# Major Edge TPU Challenges

We find that the accelerator suffers from three major challenges:

- 1 Operates **significantly below** its peak **throughput**
- 2 Operates **significantly below** its peak **energy efficiency**
- 3 Handles **memory accesses inefficiently**

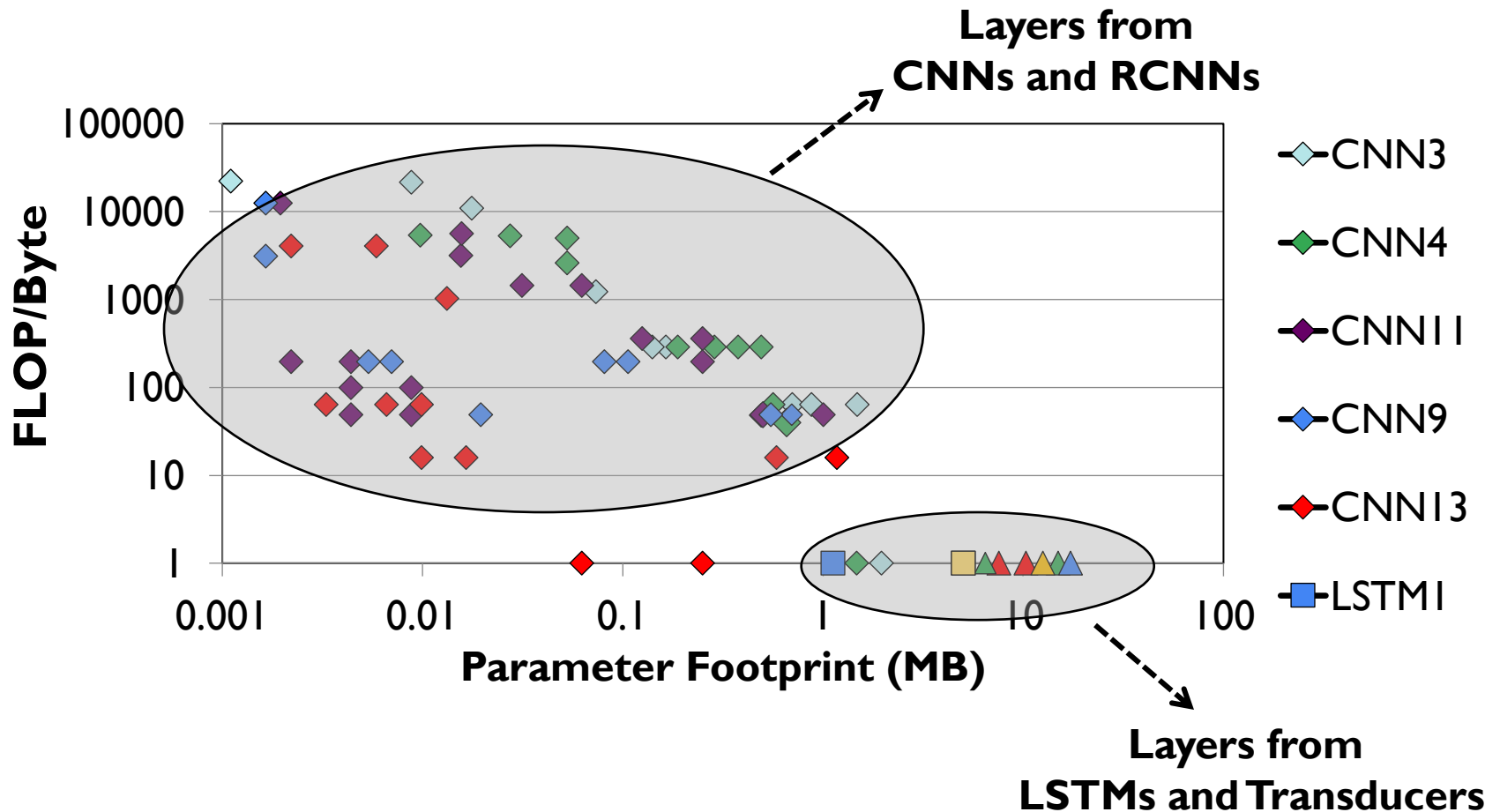
**Question: Where do these challenges come from?**



# Model Analysis: Let's Take a Deeper Look Into the Google Edge NN Models

# Diversity Across the Models

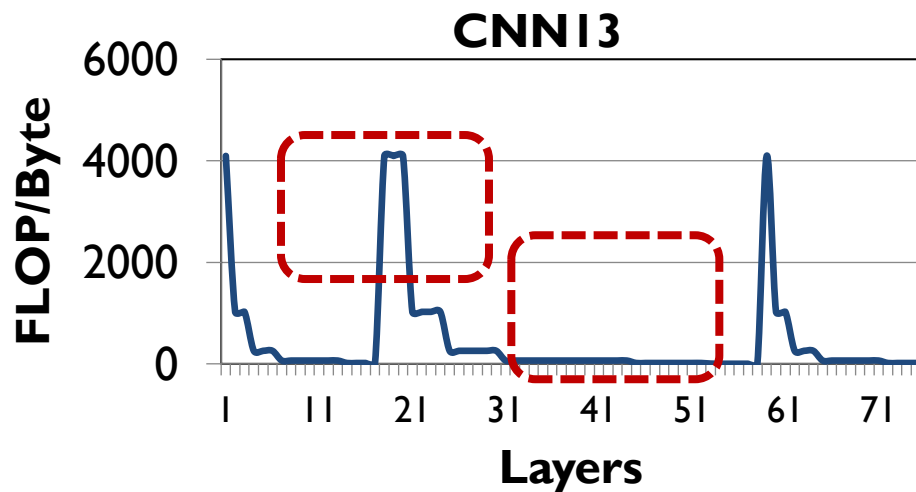
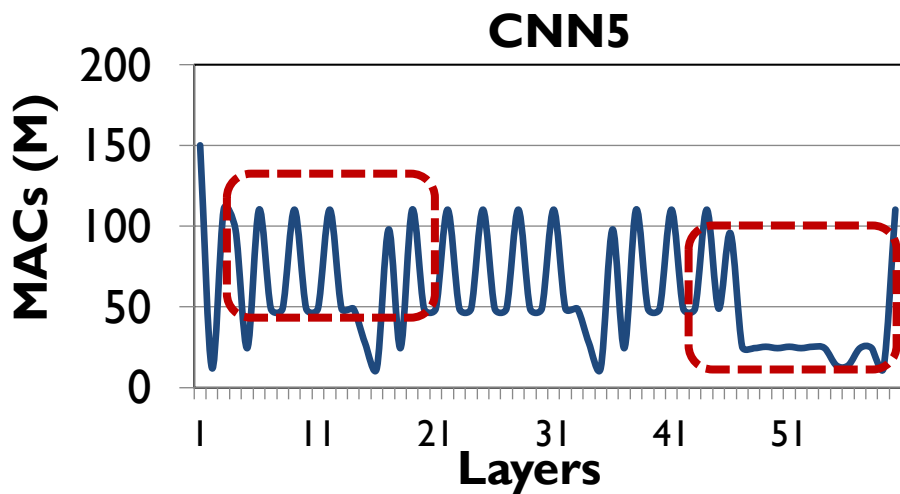
**Insight I:** there is **significant variation** in terms of layer characteristics **across the models**



# Diversity Within the Models

**Insight 2:** even **within** each model, layers exhibit **significant variation** in terms of layer characteristics

For example, our analysis of edge **CNN** models shows:

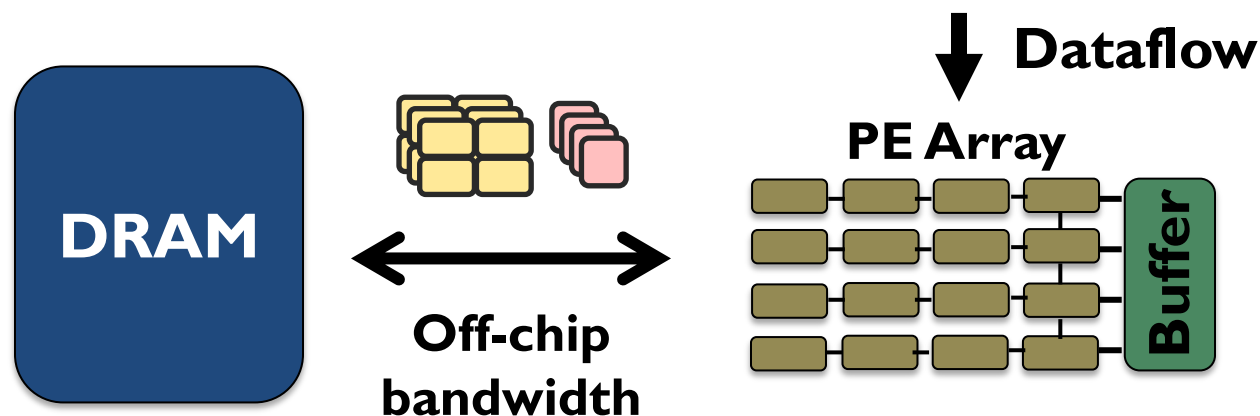


Variation in **MAC intensity**: up to **200x** across layers

Variation in **FLOP/Byte**: up to **244x** across layers

# Root Cause of Accelerator Challenges

The **key components** of Google Edge TPU are completely **oblivious** to **layer heterogeneity**



Edge accelerators typically take **a monolithic** approach: equip the accelerator with **an over-provisioned PE array** and on-chip buffer, **a rigid dataflow**, and **fixed off-chip bandwidth**



**While this approach might work for a specific group of layers, it fails to efficiently execute inference across a wide variety of edge models**

# Outline

1 Introduction

2 Edge TPU and Model Characterization

3 **Mensa Framework**

4 Mensa-G: Mensa for Google Edge Models

5 Evaluation

6 Conclusion

# Mensa Framework

**Goal:** design an edge accelerator that can efficiently run inference across **a wide range of different models** and **layers**

Instead of running the entire NN model on  
**a monolithic accelerator:**

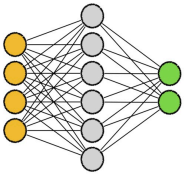


**Mensa: a new acceleration framework for edge NN inference**

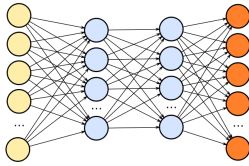
# Mensa High-Level Overview

## Edge TPU Accelerator

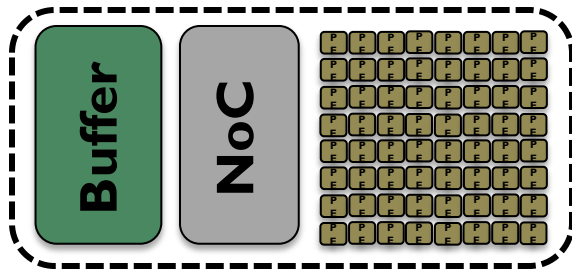
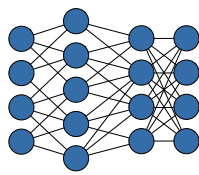
Model A



Model B



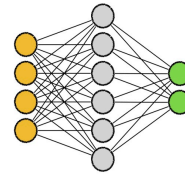
Model C



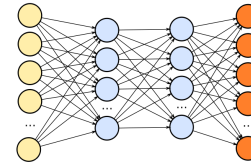
Monolithic Accelerator

## Mensa

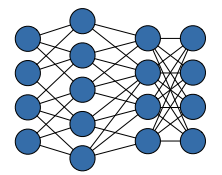
Model A



Model B



Model C

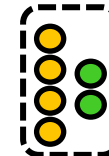


Runtime

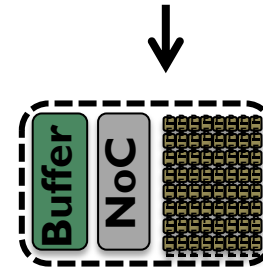
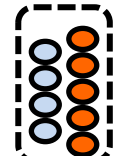
Family 1



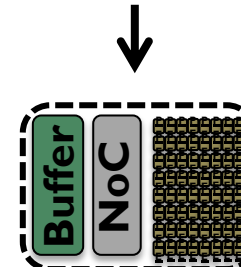
Family 2



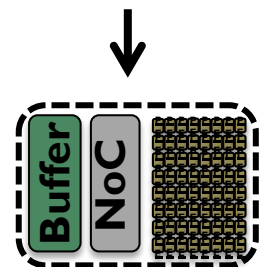
Family 3



Acc. 1



Acc. 2

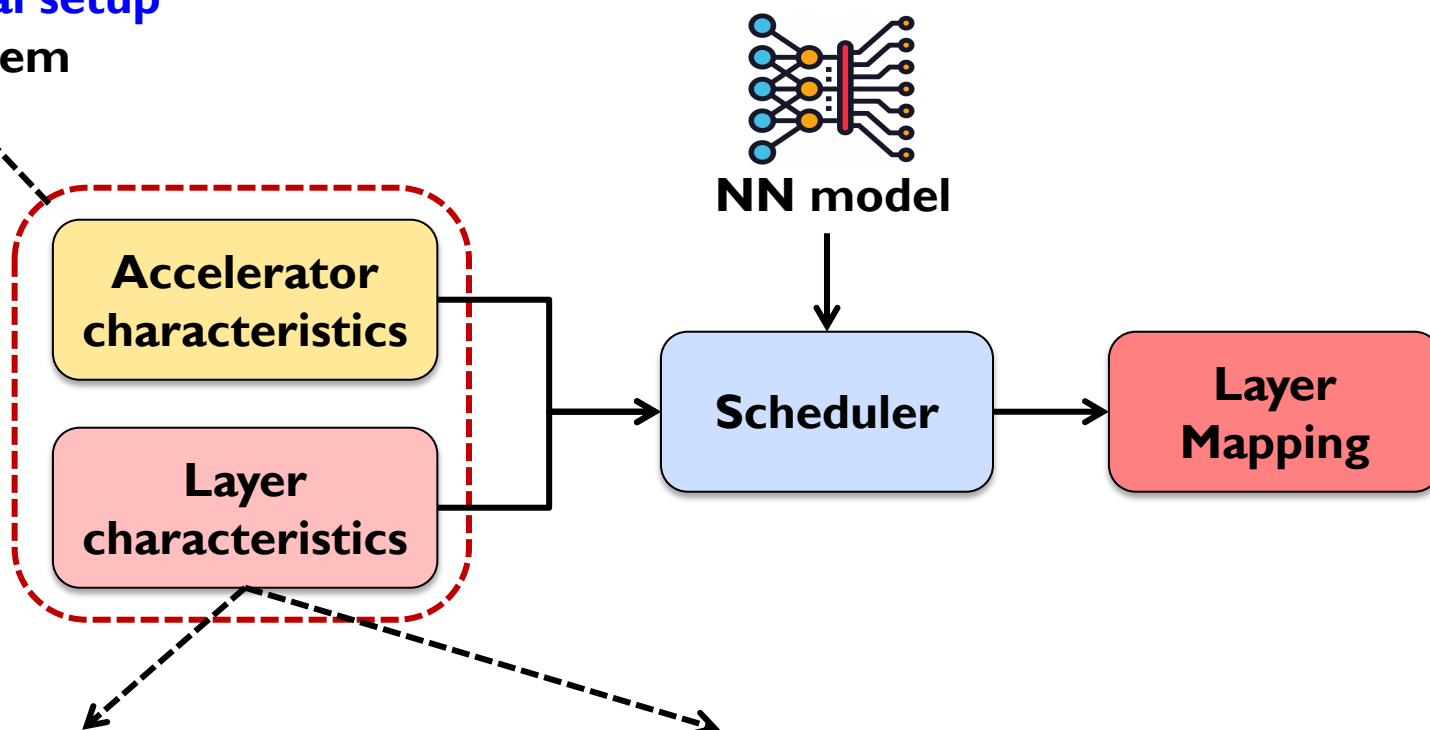


Acc. 3

# Mensa Runtime Scheduler

The **goal** of Mensa's software **runtime scheduler** is to **identify** **which accelerator** each **layer** in an NN model should run on

Generated **once**  
during **initial setup**  
of a system



Each of the accelerators  
caters to  
**a specific family** of layers

Layers tend to **group**  
together into a small  
number of **families**



# Mensa Runtime Scheduler

The **goal** of Mensa's software **runtime scheduler** is to **identify** which accelerator each **layer** in an NN model should run on

Generated **once**  
during **initial setup**

## Google Neural Network Models for Edge Devices: Analyzing and Mitigating Machine Learning Inference Bottlenecks

Amirali Boroumand<sup>†◇</sup>

Geraldo F. Oliveira<sup>\*</sup>

Saugata Ghose<sup>‡</sup>

Xiaoyu Ma<sup>§</sup>

Berkin Akin<sup>§</sup>

Eric Shiu<sup>§</sup>

Ravi Narayanaswami<sup>§</sup>

Onur Mutlu<sup>\*†</sup>

<sup>†</sup>*Carnegie Mellon Univ.*

<sup>◇</sup>*Stanford Univ.*

<sup>‡</sup>*Univ. of Illinois Urbana-Champaign*

<sup>§</sup>*Google*

<sup>\*</sup>*ETH Zürich*

Layer  
characteristics

Each of the accelerators  
caters to  
**a specific family** of layers

Layers tend to **group**  
together into a small  
number of **families**

# Outline

1 Introduction

2 Edge TPU and Model Characterization

3 Mensa Framework

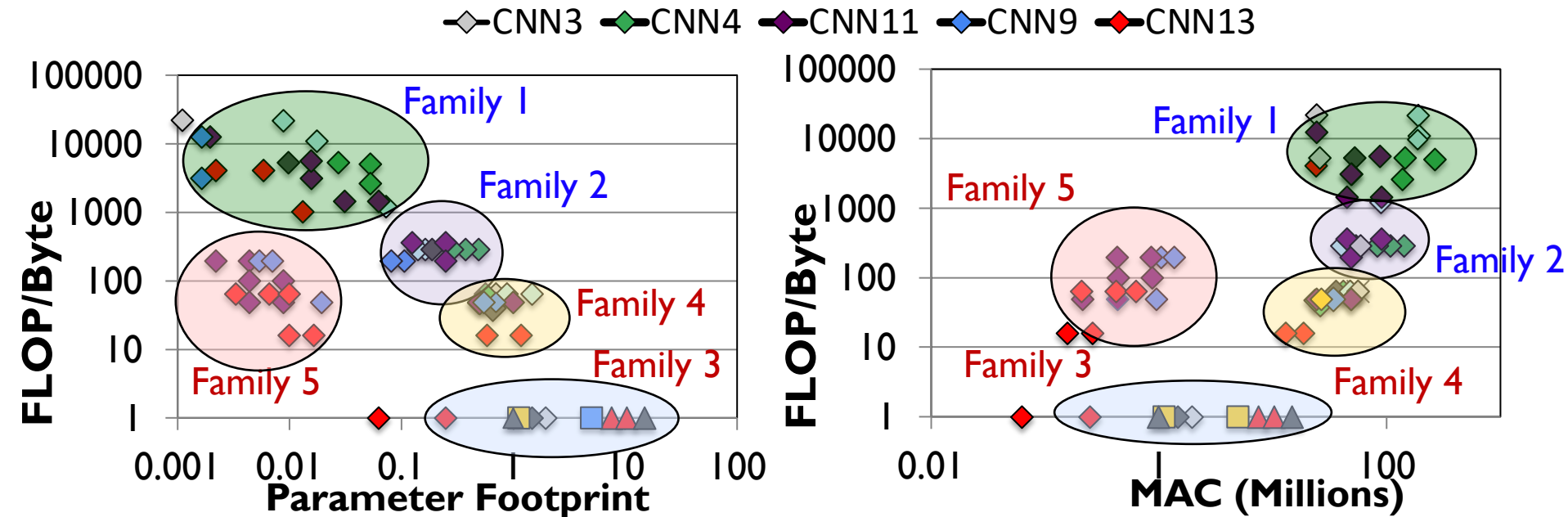
**4 Mensa-G: Mensa for Google Edge Models**

5 Evaluation

6 Conclusion

# Identifying Layer Families

**Key observation: the majority of layers group into a small number of layer families**



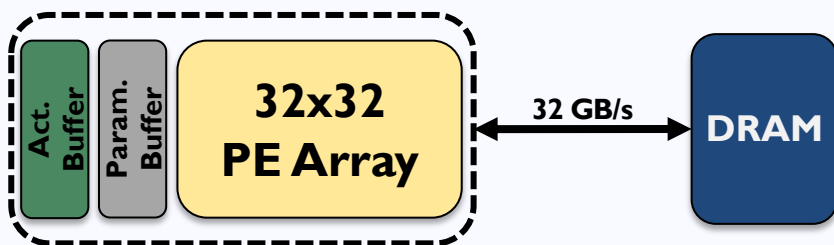
**Families 1 & 2: low parameter footprint, high data reuse and MAC intensity**  
→ compute-centric layers

**Families 3, 4 & 5: high parameter footprint, low data reuse and MAC intensity**  
→ data-centric layers

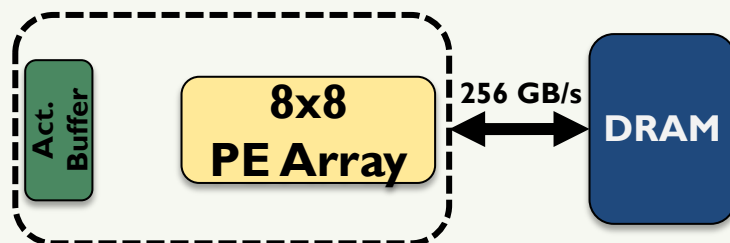
# Mensa-G: Mensa for Google Edge Models

Based on **key characteristics** of families, we design **three accelerators** to efficiently execute inference across our Google NN models

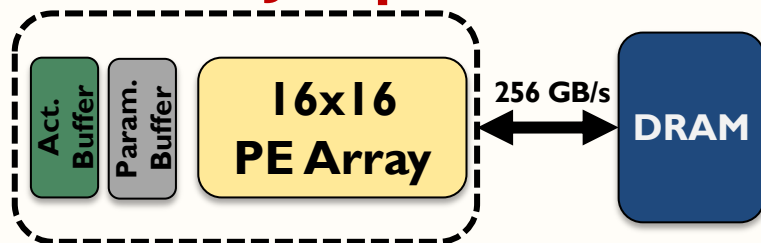
## Pascal



## Pavlov



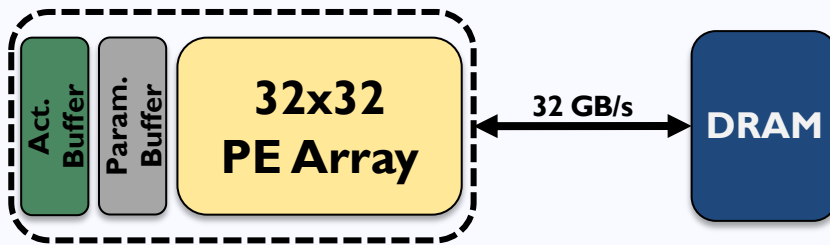
## Jacquard



# Mensa-G: Mensa for Google Edge Models

Based on **key characteristics** of families, we design **three accelerators** to efficiently execute inference across our Google NN models

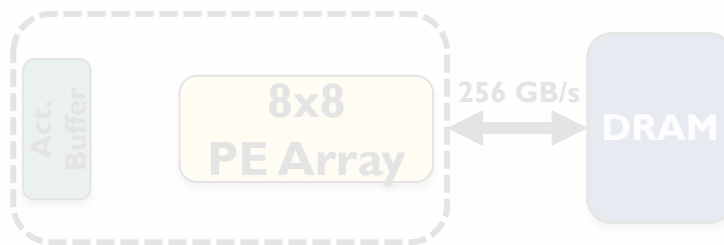
## Pascal



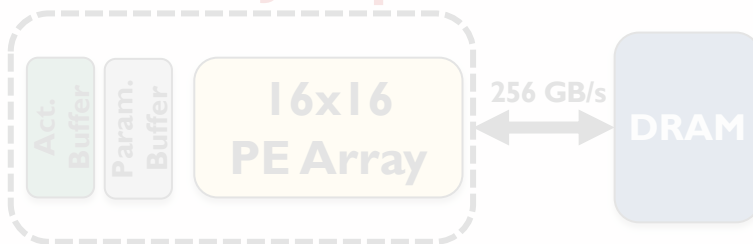
Families 1&2 → **compute-centric** layers

- **32x32 PE Array** → 2 TFLOP/s
- **256KB Act. Buffer** → **8x** Reduction
- **128KB Param. Buffer** → **32x** Reduction
- **On-chip accelerator**

## Pavlov



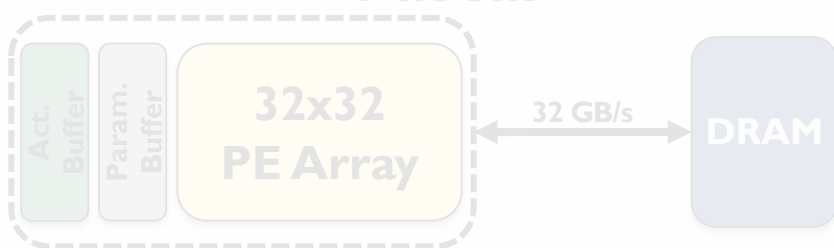
## Jacquard



# Mensa-G: Mensa for Google Edge Models

Based on **key characteristics** of families, we design **three accelerators** to efficiently execute inference across our Google NN models

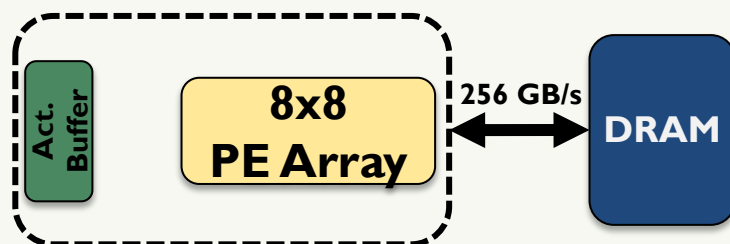
## Pascal



Families 1&2 → **compute-centric** layers

- **32x32 PE Array** → 2 TFLOP/s
- **256KB Act. Buffer** → **8x** Reduction
- **128KB Param. Buffer** → **32x** Reduction
- **On-chip accelerator**

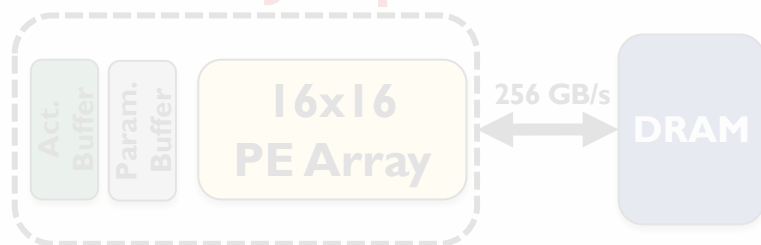
## Pavlov



Family 3 → **LSTM data-centric** layers

- **8x8 PE Array** → 128 GFLOP/s
- **128KB Act. Buffer** → **16x** Reduction
- **No** Param. Buffer → **4MB in Baseline**
- **Near-data accelerator**

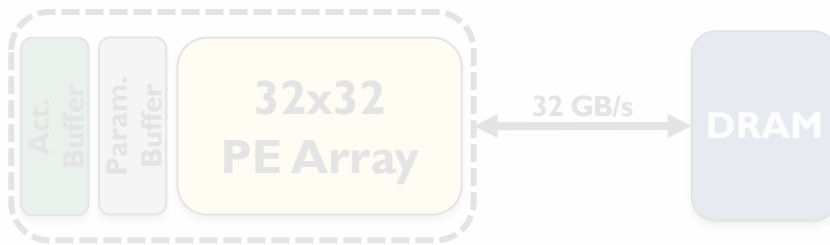
## Jacquard



# Mensa-G: Mensa for Google Edge Models

Based on **key characteristics** of families, we design **three accelerators** to efficiently execute inference across our Google NN models

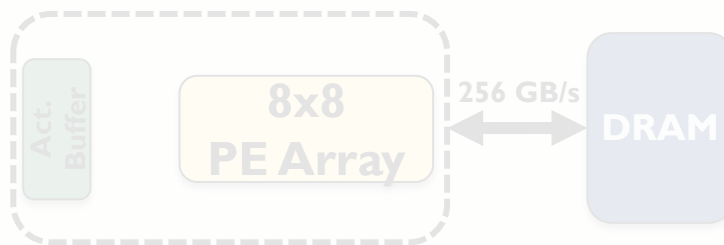
## Pascal



Families 1&2 → **compute-centric** layers

- **32x32 PE Array** → 2 TFLOP/s
- **256KB Act. Buffer** → **8x** Reduction
- **128KB Param. Buffer** → **32x** Reduction
- **On-chip accelerator**

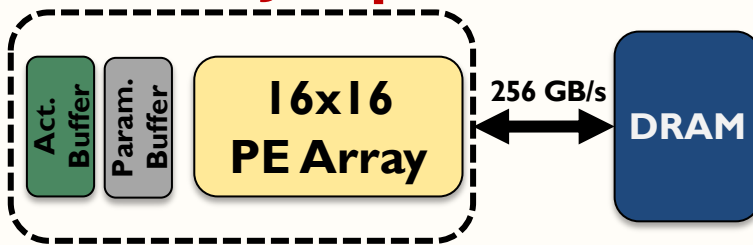
## Pavlov



Family 3 → **LSTM data-centric** layers

- **8x8 PE Array** → 128 GFLOP/s
- **128KB Act. Buffer** → **16x** Reduction
- **No** Param. Buffer → **4MB in Baseline**
- **Near-data accelerator**

## Jacquard



Families 4&5 → **non-LSTM data-centric** layers

- **16x16 PE Array** → 256 GFLOP/s
- **128KB Act. Buffer** → **16x** Reduction
- **128KB Param. Buffer** → **32x** Reduction
- **Near-data accelerator**

# Mensa-G: Mensa for Google Edge Models

Based on **key characteristics** of families, we design **three accelerators** to efficiently execute inference across our Google NN models

## Pascal

Families 1&2 → **compute-centric** layers

- **32x32 PE Array** → 2 TFLOP/s

- **256KB Act. Buffer** → **8x** Reduction

## Google Neural Network Models for Edge Devices: Analyzing and Mitigating Machine Learning Inference Bottlenecks

Amirali Boroumand<sup>†◇</sup>

Saugata Ghose<sup>‡</sup>

Berkin Akin<sup>§</sup>

Ravi Narayanaswami<sup>§</sup>

Geraldo F. Oliveira<sup>\*</sup>

Xiaoyu Ma<sup>§</sup>

Eric Shiu<sup>§</sup>

Onur Mutlu<sup>\*†</sup>

<sup>†</sup>Carnegie Mellon Univ.

<sup>◇</sup>Stanford Univ.

<sup>‡</sup>Univ. of Illinois Urbana-Champaign

<sup>§</sup>Google

<sup>\*</sup>ETH Zürich

- **Near-data accelerator**

## Jacquard

Families 4&5 → **non-LSTM data-centric** layers

- **16x16 PE Array** → 256 GFLOP/s

- **128KB Act. Buffer** → **16x** Reduction

- **128KB Param. Buffer** → **32x** Reduction

- **Near-data accelerator**





# Outline

1 Introduction

2 Edge TPU and Model Characterization

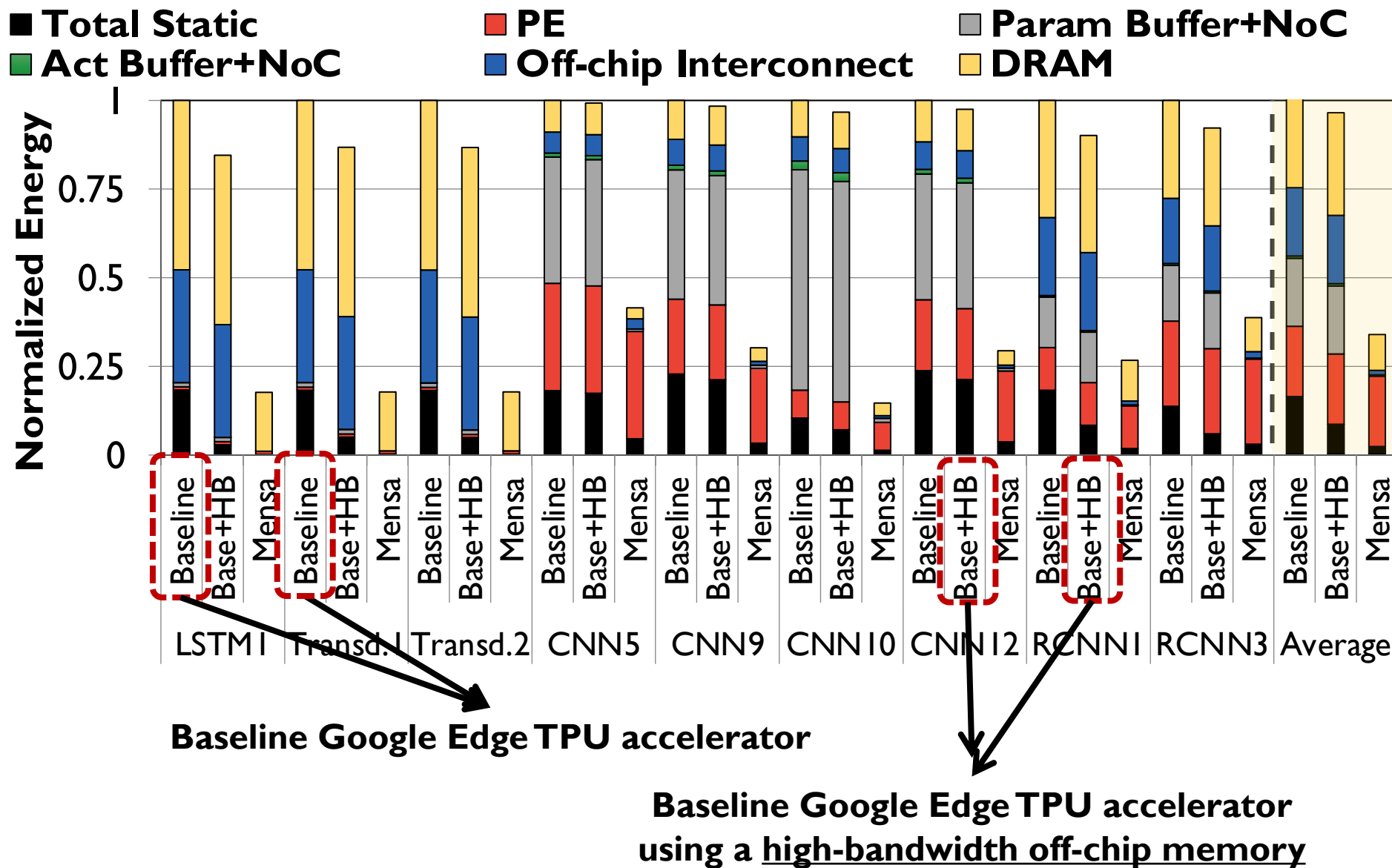
3 Mensa Framework

4 Mensa-G: Mensa for Google Edge Models

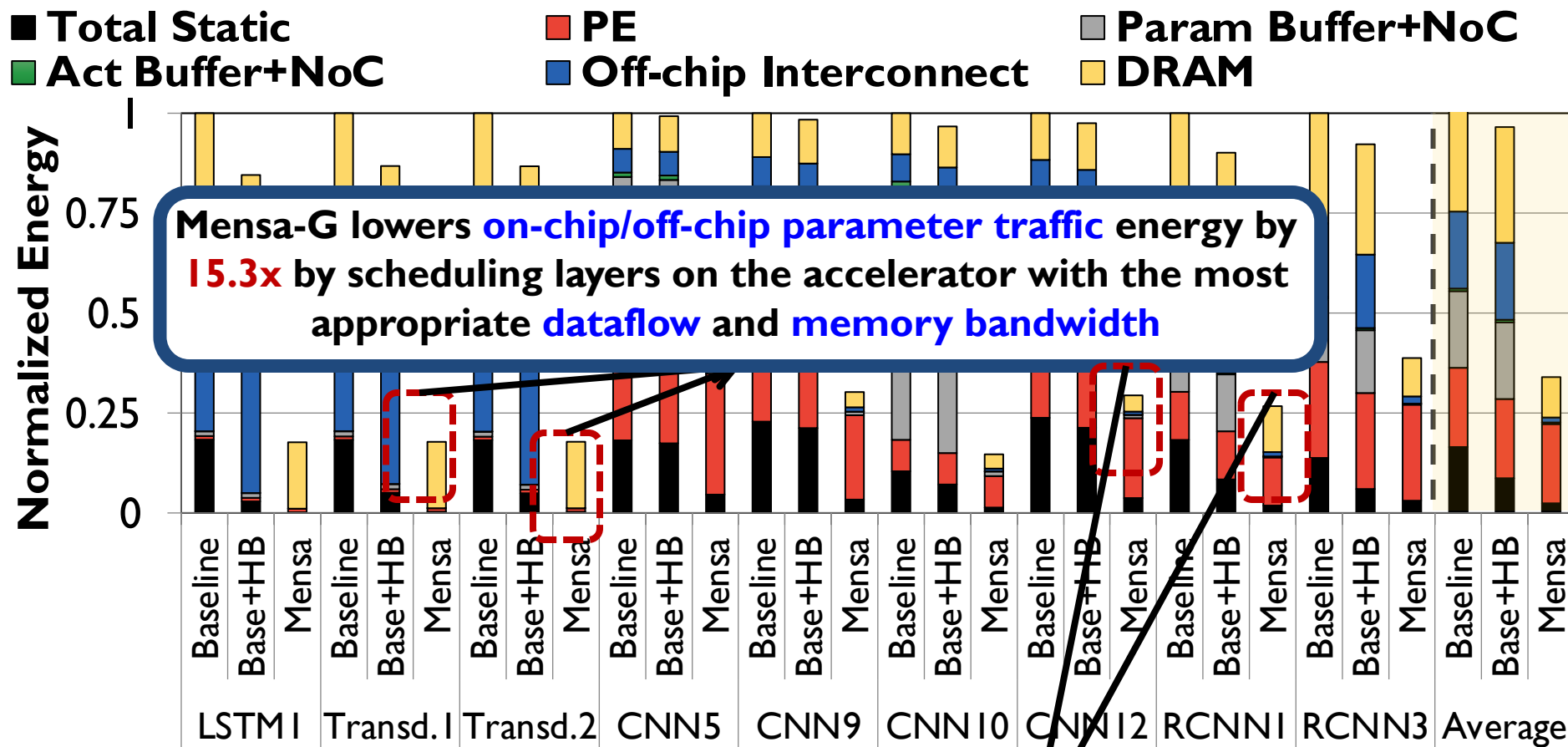
5 Evaluation

6 Conclusion

# Energy Analysis

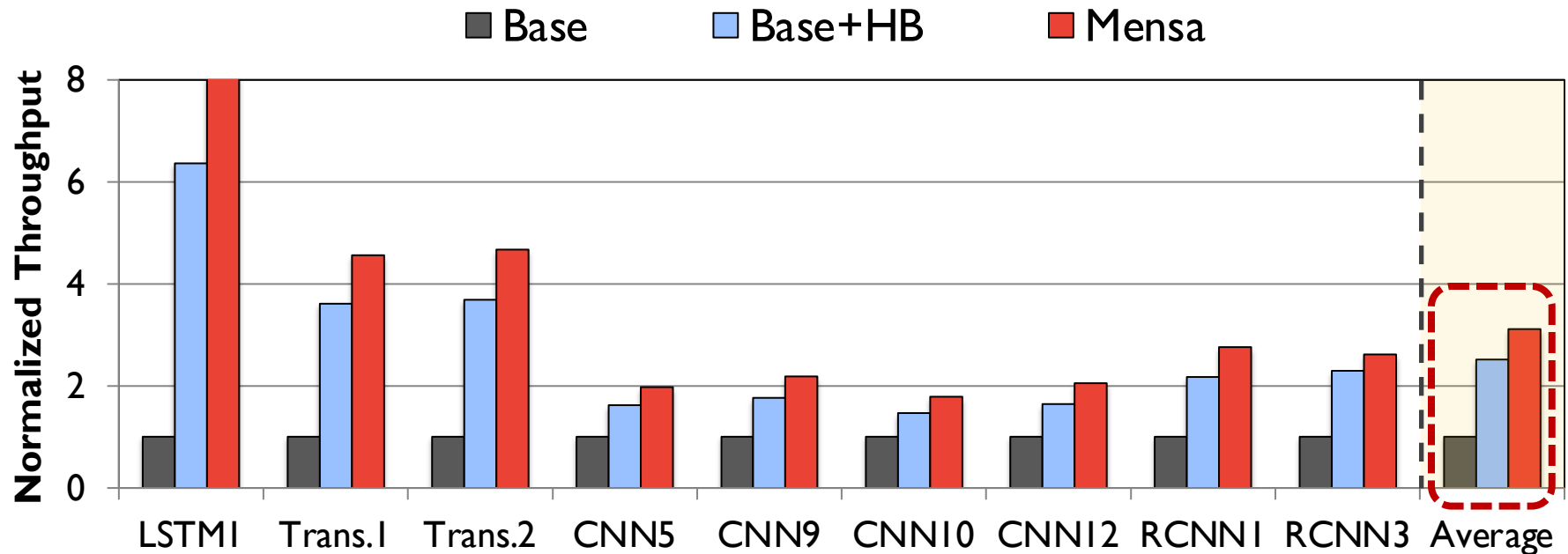


# Energy Analysis



Mensa-G improves energy efficiency by 3.0X compared to the Baseline

# Throughput Analysis



**Mensa-G improves throughput by 3.1X compared to the Baseline**

# More in the Paper

- **Details about Mensa Runtime Scheduler**
- **Details about Pascal, Pavlov, and Jacquard's dataflows**
- **Energy comparison with Eyeriss v2**
- **Mensa-G's utilization results**
- **Mensa-G's inference latency results**

# More in the Paper

- Details about Mensa Runtime Scheduler

- Details about Pascal Boyer and Jeannot

## Google Neural Network Models for Edge Devices: Analyzing and Mitigating Machine Learning Inference Bottlenecks

Amirali Boroumand<sup>†◇</sup>

Saugata Ghose<sup>‡</sup>

Berkin Akin<sup>§</sup>

Ravi Narayanaswami<sup>§</sup>

Geraldo F. Oliveira<sup>\*</sup>

Xiaoyu Ma<sup>§</sup>

Eric Shiu<sup>§</sup>

Onur Mutlu<sup>\*†</sup>

<sup>†</sup>*Carnegie Mellon Univ.*

<sup>◇</sup>*Stanford Univ.*

<sup>‡</sup>*Univ. of Illinois Urbana-Champaign*

<sup>§</sup>*Google*

<sup>\*</sup>*ETH Zürich*

- Mensa-G's utilization results
- Mensa-G's inference latency results

# Outline

1 Introduction

2 Edge TPU and Model Characterization

3 Mensa Framework

4 Mensa-G: Mensa for Google Edge Models

5 Evaluation

6 Conclusion

# Conclusion

**Context:** We extensively analyze a state-of-the-art edge ML accelerator (Google Edge TPU) using 24 Google edge models

- Wide range of models (CNNs, LSTMs, Transducers, RCNNs)

**Problem:** The Edge TPU accelerator suffers from **three challenges:**

- It operates **significantly below** its peak throughput
- It operates **significantly below** its theoretical energy efficiency
- It **inefficiently** handles memory accesses

**Key Insight:** These shortcomings arise from **the monolithic design** of the Edge TPU accelerator

- The Edge TPU accelerator design does not account for **layer heterogeneity**

**Key Mechanism:** A new framework called **Mensa**

- Mensa consists of heterogeneous accelerators whose dataflow and hardware are specialized for specific families of layers

**Key Results:** We design a version of Mensa for Google edge ML models

- Mensa improves performance and energy by **3.0X** and **3.1X**
- Mensa reduces cost and improves area efficiency



# Google Neural Network Models for Edge Devices: Analyzing and Mitigating Machine Learning Inference Bottlenecks

P&S Processing-in-Memory  
Fall 2022  
17 January 2023



**Amirali Boroumand**

**Saugata Ghose**

**Berkin Akin**

**Ravi Narayanaswami**

**Geraldo F. Oliveira**

**Xiaoyu Ma**

**Eric Shiu**

**Onur Mutlu**

**SAFARI**

# Polynesia:

## Enabling High-Performance and Energy-Efficient Hybrid Transactional/Analytical Databases with Hardware/Software Co-Design

P&S Processing-in-Memory

Fall 2022

17 January 2023

**Amirali Boroumand**  
**Geraldo F. Oliveira**

**Saugata Ghose**  
**Onur Mutlu**

# Executive Summary

- **Context:** Many applications need to perform real-time data analysis using an Hybrid Transactional/Analytical Processing (HTAP) system
  - An ideal HTAP system should have **three properties**:  
(1) **data freshness** and **consistency**, (2) **workload-specific optimization**,  
(3) **performance isolation**
- **Problem:** Prior works **cannot achieve all properties** of an ideal HTAP system
- **Key Idea:** Divide the system into transactional and analytical **processing islands**
  - Enables **workload-specific optimizations** and **performance isolation**
- **Key Mechanism:** Polynesia, a novel hardware/software cooperative design for in-memory HTAP databases
  - Implements **custom algorithms and hardware** to reduce the costs of **data freshness** and **consistency**
  - Exploits **PIM** for analytical processing to alleviate **data movement**
- **Key Results:** Polynesia outperforms three state-of-the-art HTAP systems
  - Average transactional/analytical throughput improvements of **1.7x/3.7x**
  - **48%** reduction on energy consumption

# Outline

1

Introduction

2

Limitations of HTAP Systems

3

Polynesia: Overview

4

Update Propagation Mechanism

5

Consistency Mechanism

6

Analytical Engine

7

Evaluation

8

Conclusion

# Outline

1

Introduction

2

Limitations of HTAP Systems

3

Polynesia: Overview

4

Update Propagation Mechanism

5

Consistency Mechanism

6

Analytical Engine

7

Evaluation

8

Conclusion

# Real-Time Analysis

An explosive interest in many applications domains to perform data analytics on the most recent version of data (real-time analysis)

Use **transactions** to **record** each periodic sample of data from **all sensors**

Run **analytics** across sensor data to make **real-time** steering decisions

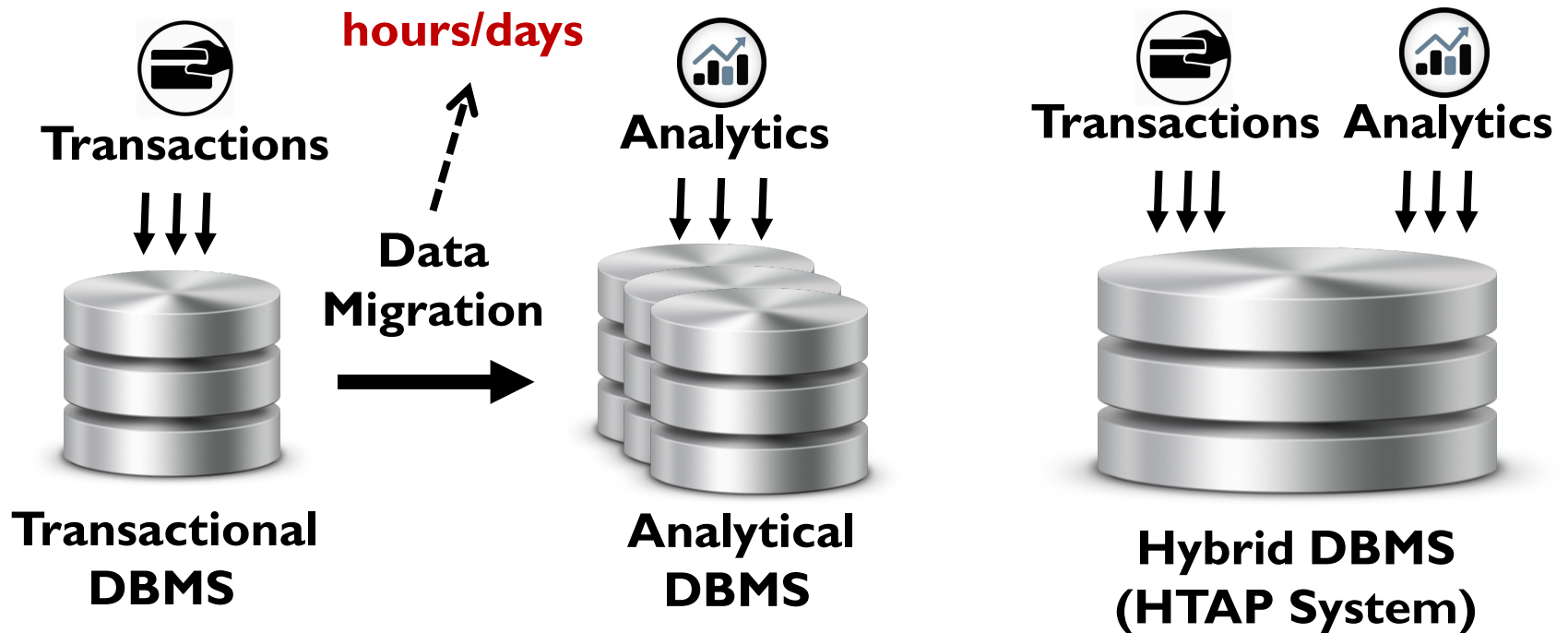


Self-Driving Cars

For these applications, it is **critical** to analyze **the transactions** in **real-time** as the data's value **diminishes** over time

# HTAP: Supporting Real-Time Analysis

Traditionally, **new transactions (updates)** are propagated to the **analytical database** using a **periodic** and **costly** process



To support real-time analysis: a single hybrid DBMS is used to execute both transactional and analytical workloads

# Ideal HTAP System Properties

An ideal HTAP system should have **three properties**:

## 1 Workload-Specific Optimizations

- Transactional and analytical workloads must benefit from their **own specific optimizations**

## 2 Data Freshness and Consistency Guarantees

- Guarantee access to the **most recent version of data** for analytics while ensuring that transactional and analytical workloads have a **consistent** view of data

## 3 Performance Isolation

- Latency and throughput of transactional and analytical workloads are the same as if they were **run in isolation**

**Achieving all three properties at the same time is very challenging**



# Outline

1

Introduction

2

**Limitations of HTAP Systems**

3

Polynesia: Overview

4

Update Propagation Mechanism

5

Consistency Mechanism

6

Analytical Engine

7

Evaluation

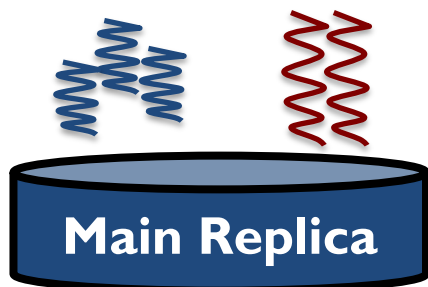
8

Conclusion

# State-of-the-Art HTAP Systems

We study two major types of HTAP systems:

Transactions Analytics

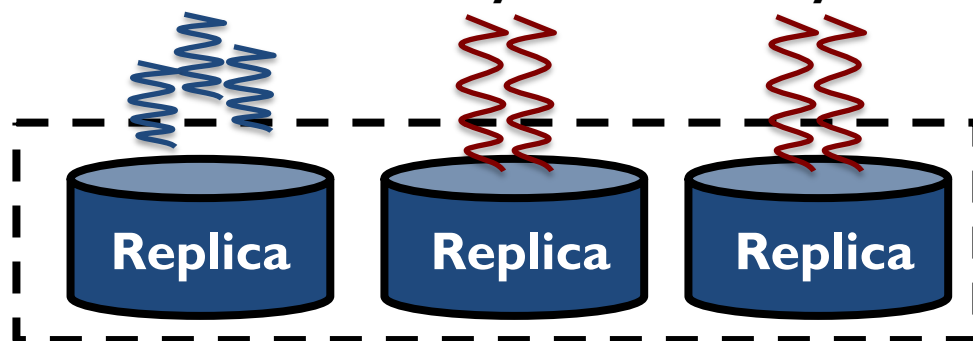


Single-Instance

Transactions

Analytics

Analytics



Multiple-Instance

We observe **two key problems**:

1

Data freshness and consistency mechanisms are costly and cause a drastic reduction in throughput

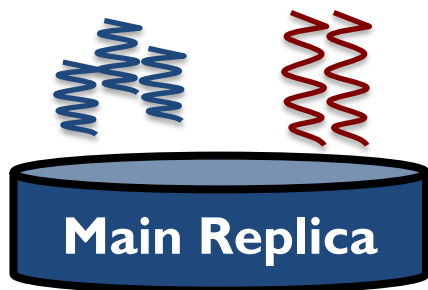
2

These systems fail to provide performance isolation because of high main memory contention

# State-of-the-Art HTAP Systems

We study two major types of HTAP systems:

Transactions Analytics

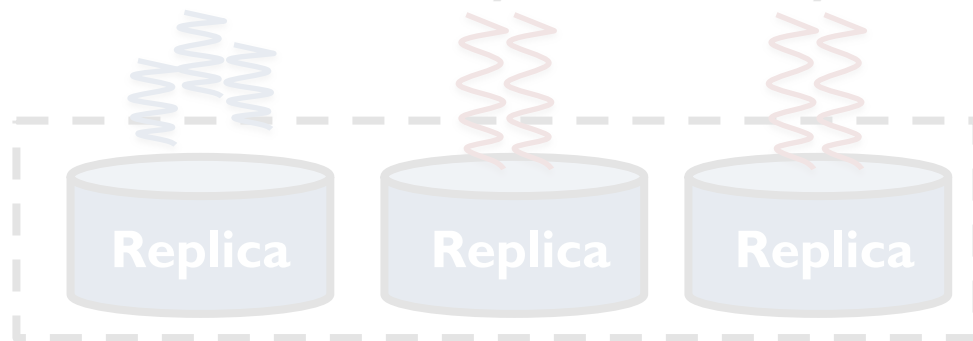


Single-Instance

Transactions

Analytics

Analytics



Multiple-Instance

We observe **two key problems**:

1

Data freshness and consistency mechanisms are costly and cause a drastic reduction in throughput

2

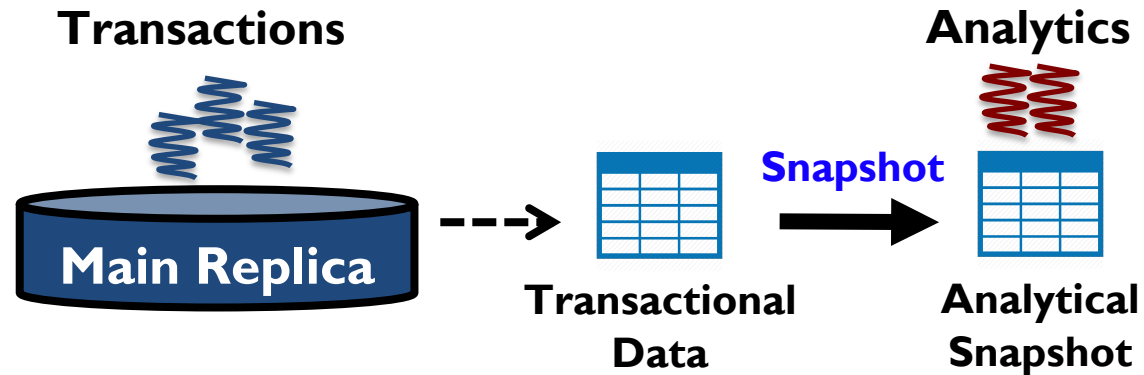
These systems fail to provide performance isolation because of high main memory contention

# Single-Instance: Data Consistency

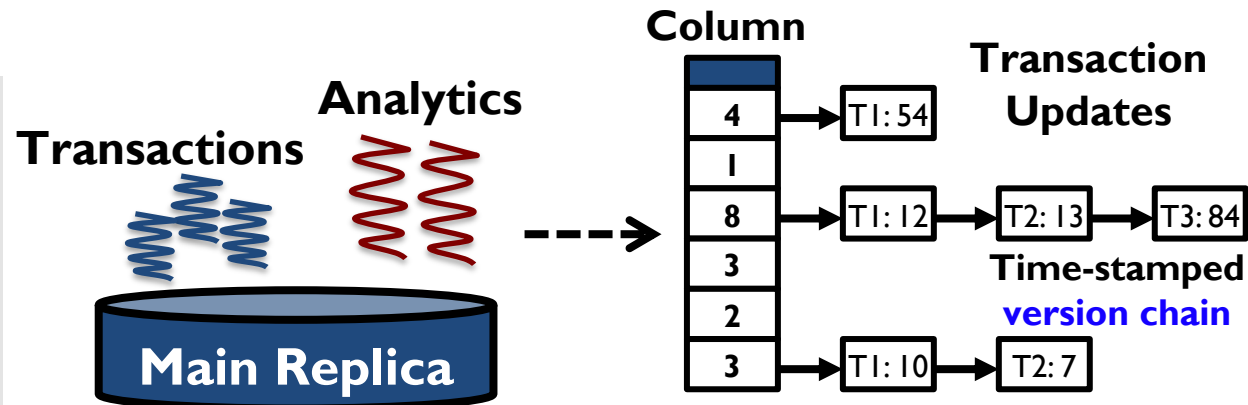
Since both **analytics** and **transactions** work on the **same data concurrently**, we need to ensure that the data is **consistent**

There are **two major mechanisms** to ensure consistency:

## 1 Snapshotting

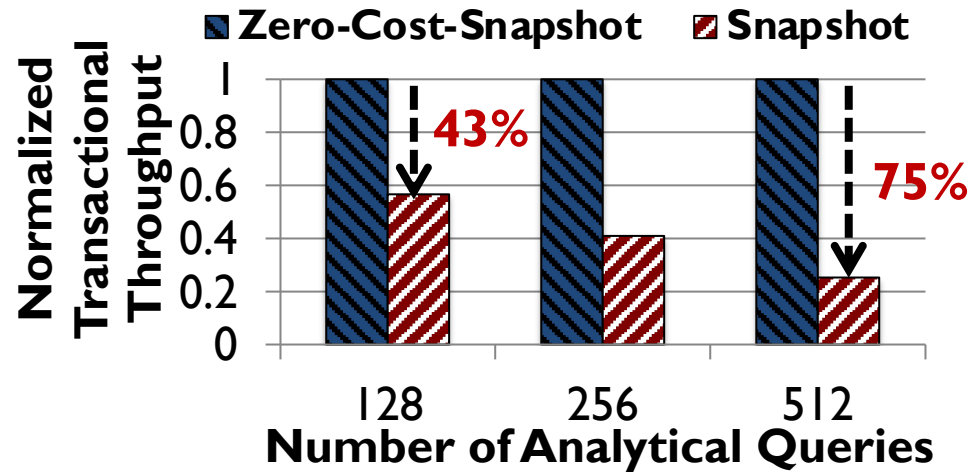


## 2 Multi-Version Concurrency Control (MVCC)

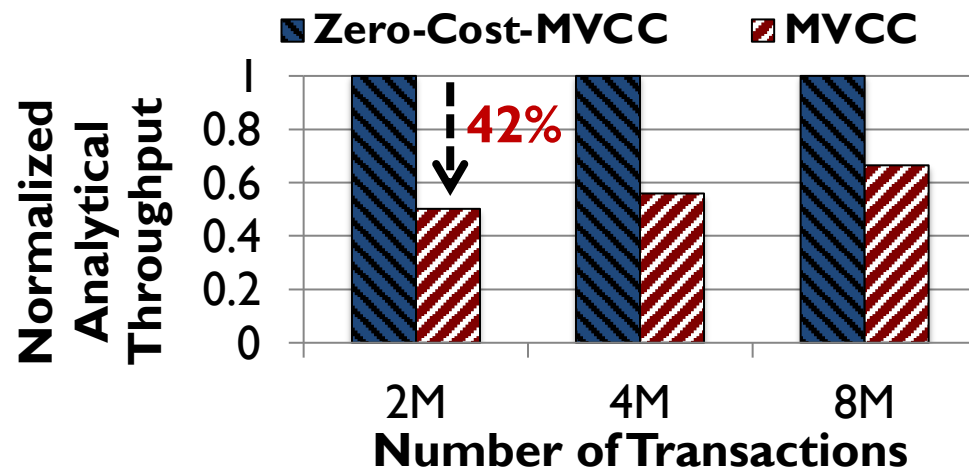


# Drawbacks of Snapshotting and MVCC

We evaluate the **throughput loss** caused by Snapshotting and MVCC:



Throughput loss comes from memcpy operation:  
generates a large amount of data movement

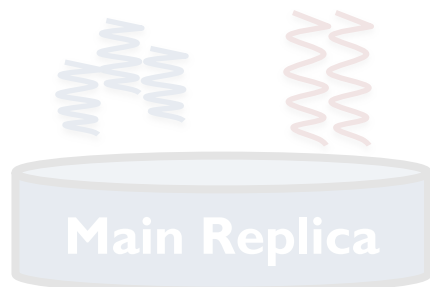


Throughput loss comes from long version chains:  
expensive time-stamp comparison and  
a large number of random memory accesses

# State-of-the-Art HTAP Systems

We study two major types of HTAP systems:

Transactions Analytics

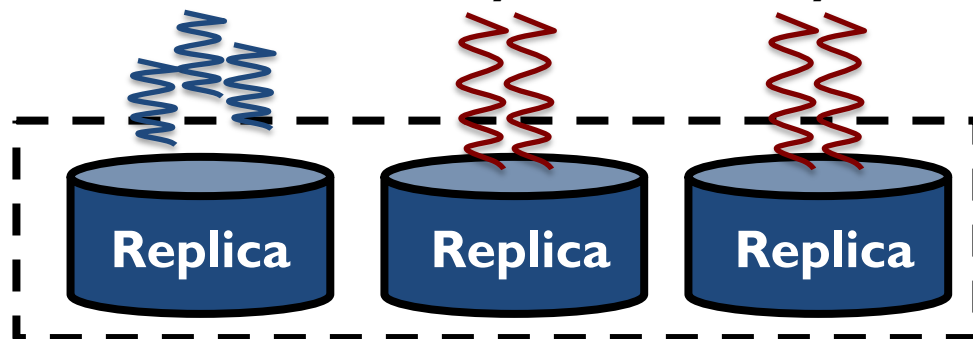


Single-Instance

Transactions

Analytics

Analytics



Multiple-Instance

We observe **two key problems**:

1

Data freshness and consistency mechanisms are costly and cause a drastic reduction in throughput

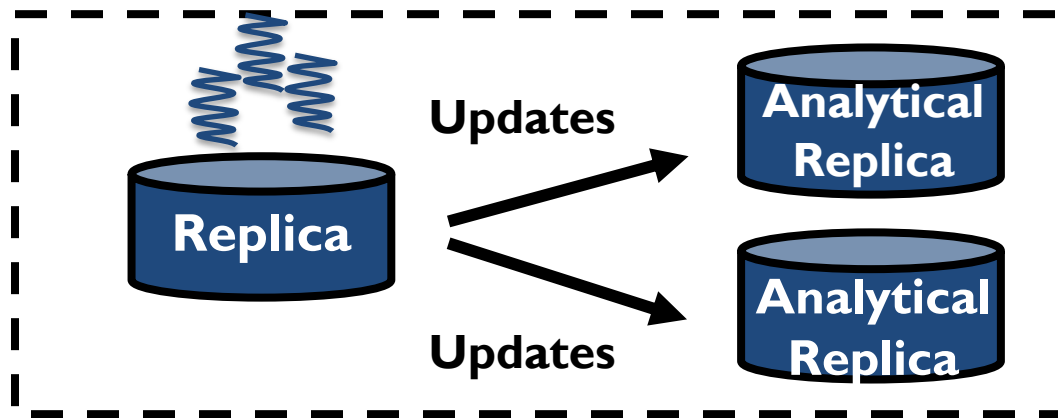
2

These systems fail to provide performance isolation because of high main memory contention

# Maintaining Data Freshness

One of the **major challenges** in multiple-instance systems is to keep **analytical** replicas **up-to-date**

Transactional queries



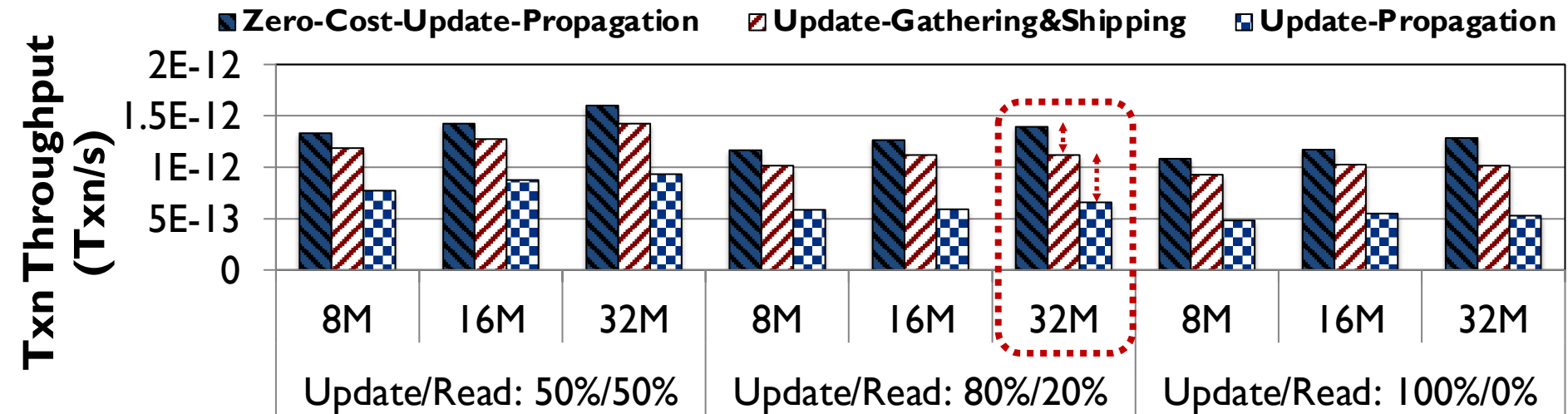
Multiple-Instance HTAP System

To maintain data freshness (via **Update Propagation**):

- 1 **Update Gathering and Shipping:** gather updates from transactional threads and ship them to analytical the replica
- 2 **Update Application:** perform the necessary format conversation and apply those updates to analytical replicas

# Cost of Update Propagation

We evaluate the **throughput loss** caused by Update Propagation:



Transactional throughput reduces by up to 21.2% during the update gathering & shipping process

Transactional throughput reduces by up to 64.2% during the update application process



# Problem and Goal

## Problems:

- 1 State-of-the-art HTAP systems **do not** achieve all of the desired HTAP properties
- 2 Data freshness and consistency mechanisms are **data-intensive** and cause a drastic **reduction** in throughput
- 3 These systems **fail** to provide **performance isolation** because of **high main memory contention**

## Goal:

Take advantage of **custom algorithm** and **processing-in-memory (PIM)** to address these **challenges**

# Outline

1

Introduction

2

Limitations of HTAP Systems

3

**Polynesia: Overview**

4

Update Propagation Mechanism

5

Consistency Mechanism

6

Analytical Engine

7

Evaluation

8

Conclusion

# Polynesia

**Key idea:** **partition** computing resources into two types of **isolated** and **specialized processing islands**

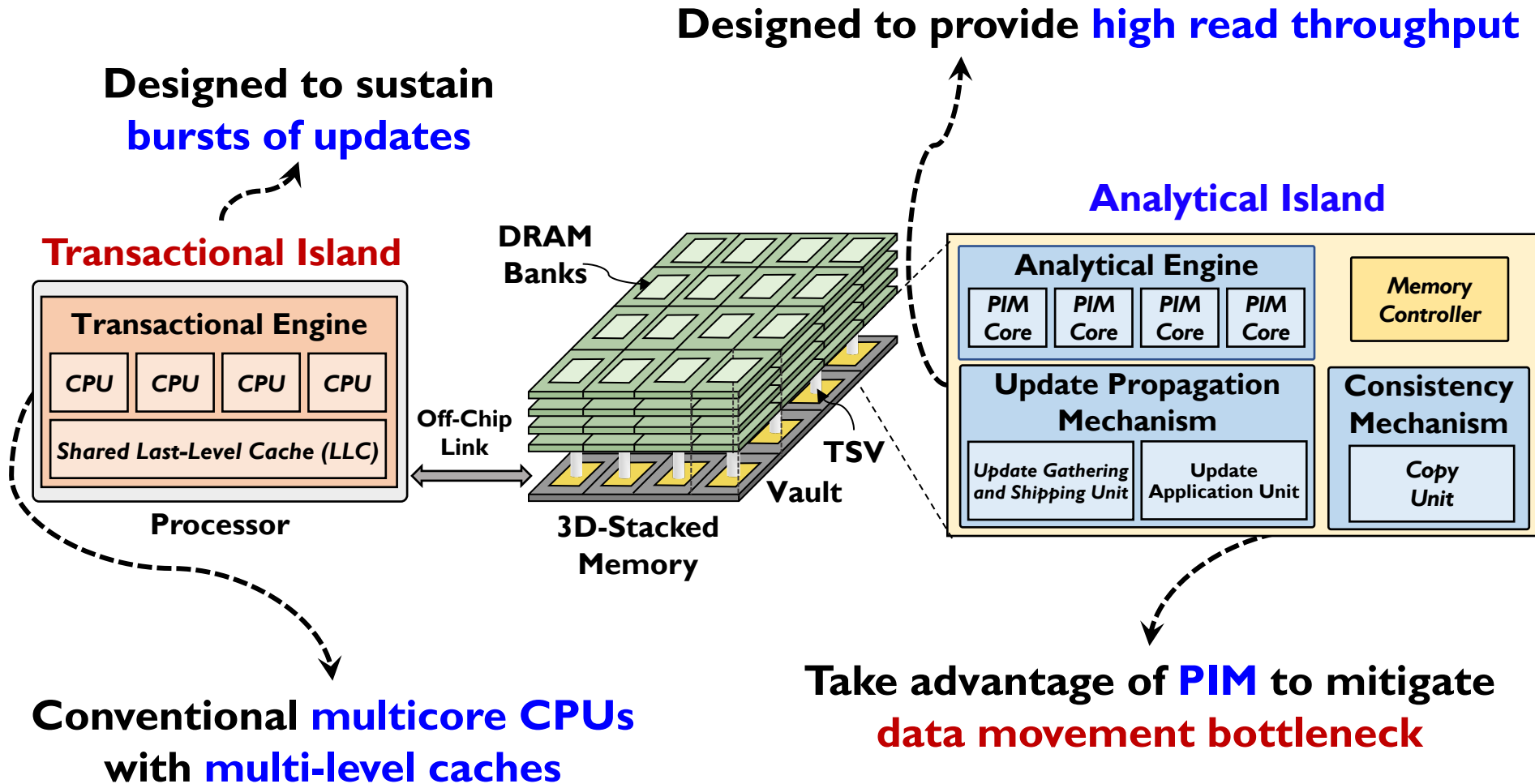


Isolating **transactional islands** from **analytical islands** allows us to:

- 1 Apply **workload-specific optimizations** to each island
- 2 Avoid high **main memory contention**
- 3 Design efficient **data freshness and consistency mechanisms** without incurring **high data movement costs**
  - Leverage **processing-in-memory (PIM)** to reduce **data movement**
  - **PIM** mitigates **data movement overheads** by placing **computation units nearby** or **inside memory**

# Polynesia: High-Level Overview

Each island includes (1) a **replica** of data, (2) an **optimized** execution engine, and (3) a set of **hardware resources**



# Outline

1

Introduction

2

Limitations of HTAP Systems

3

Polynesia: Overview

4

**Update Propagation Mechanism**

5

Consistency Mechanism

6

Analytical Engine

7

Evaluation

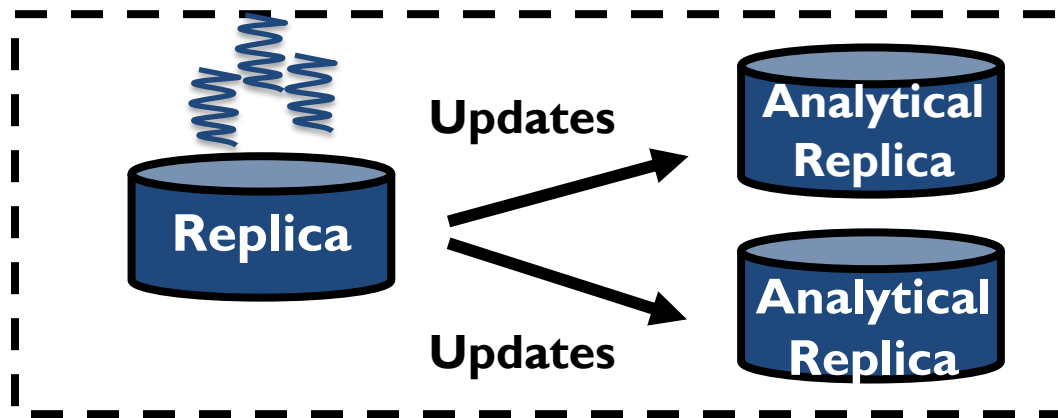
8

Conclusion

# Maintaining Data Freshness

One of the **major challenges** in multiple-instance systems is to keep **analytical** replicas **up-to-date**

Transactional queries



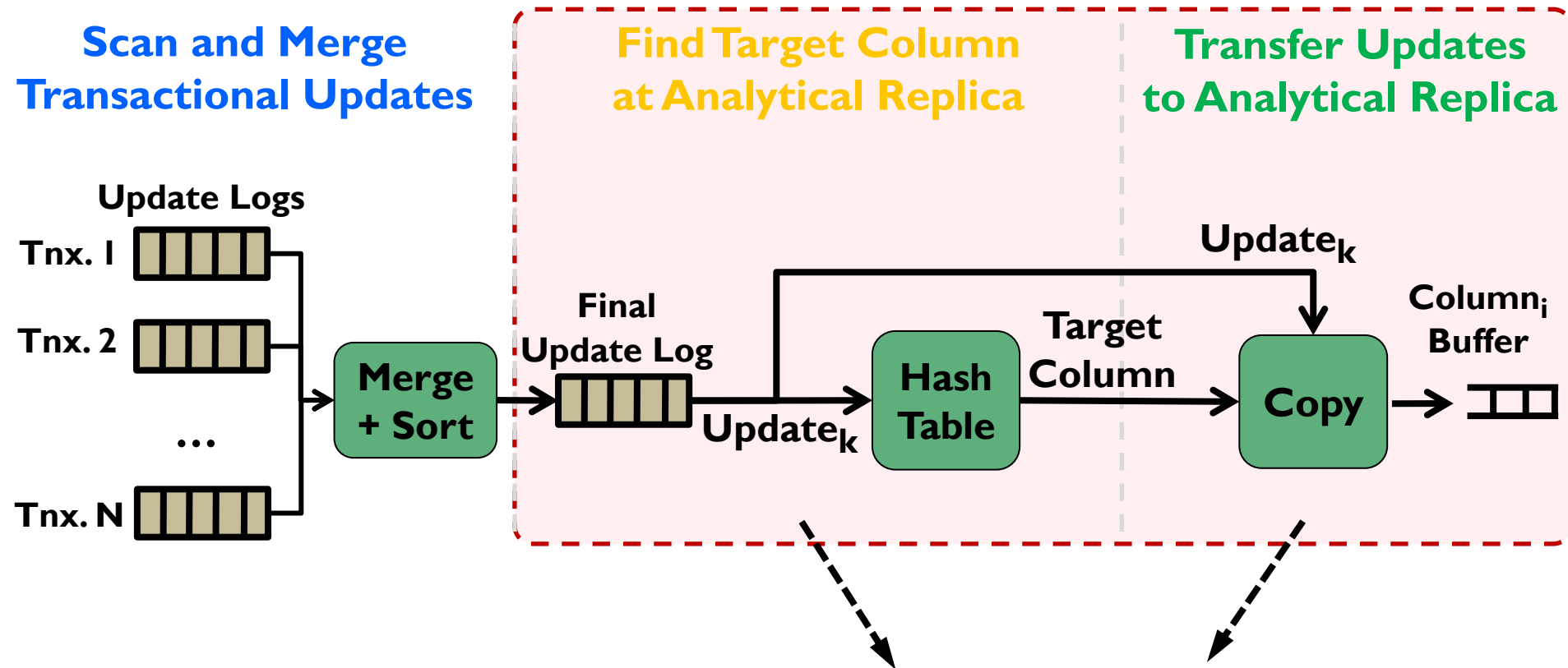
Multiple-Instance HTAP System

To maintain data freshness (via **Update Propagation**):

- 1 **Update Gathering and Shipping**: gather updates from transactional threads and ship them to analytical the replica
- 2 **Update Application**: perform the necessary format conversation and apply those updates to analytical replicas

# Update Gathering & Shipping: Algorithm

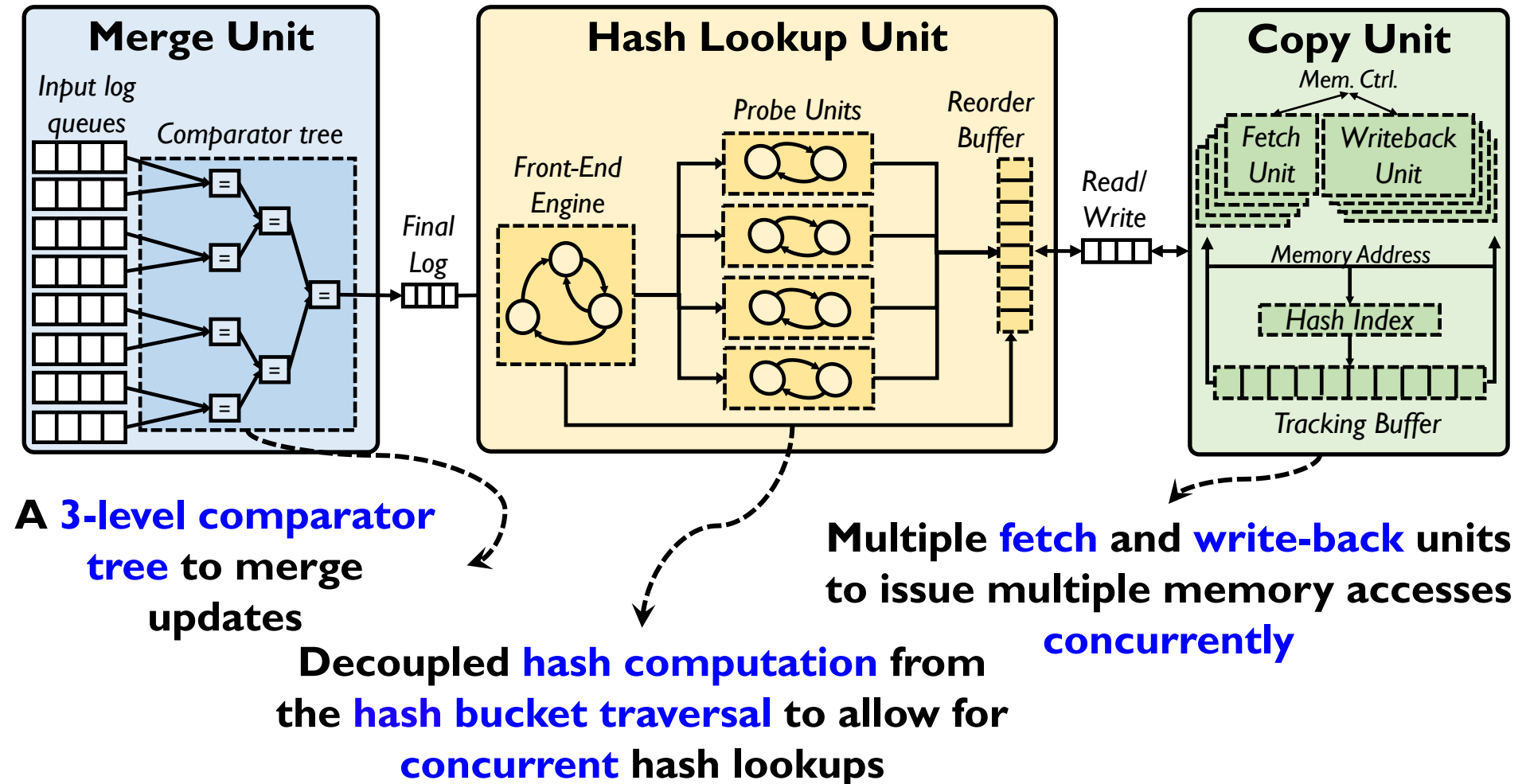
Update gathering & shipping algorithm has **three major** stages:



**2<sup>nd</sup> and 3<sup>rd</sup> stages generate a large amount of data movement and account for 87.2% of our algorithm's execution time**

# Update Gathering & Shipping: Hardware

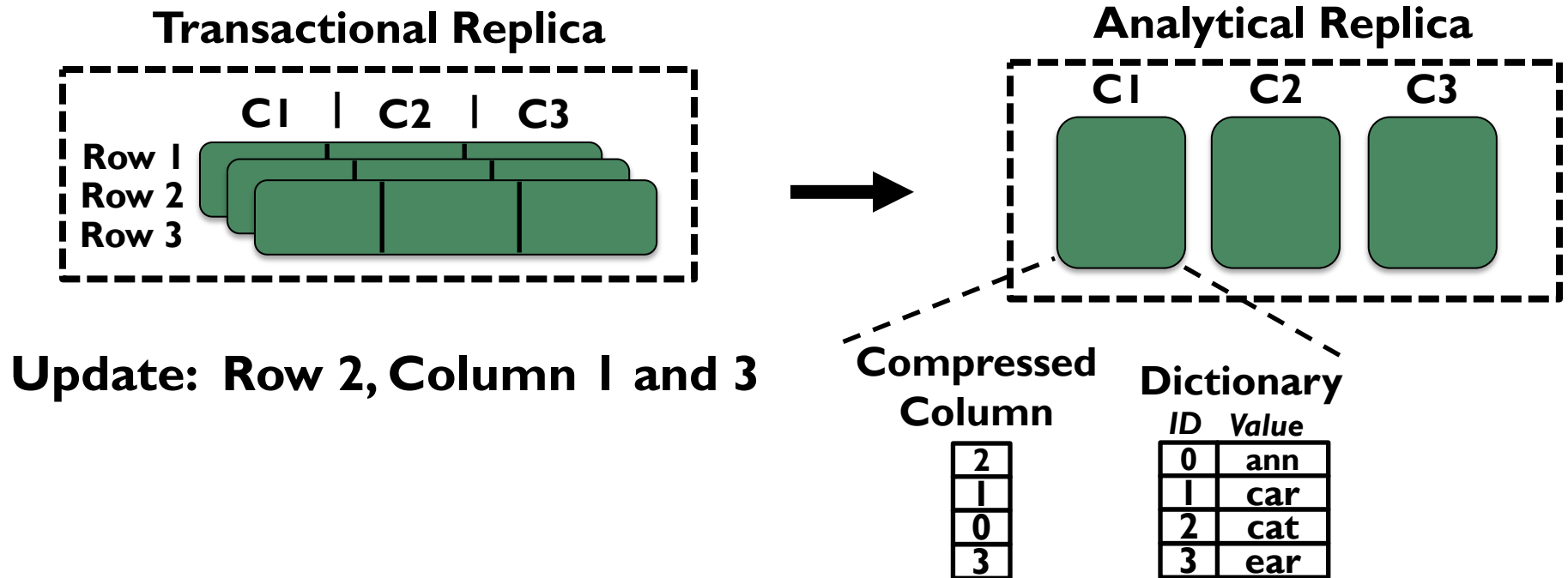
To avoid these **bottlenecks**, we design a new hardware accelerator, called **update gathering & shipping unit**





# Update Propagation: Update Application

**Goal:** perform the necessary **format conversation** and **apply** transactional updates to analytical replicas



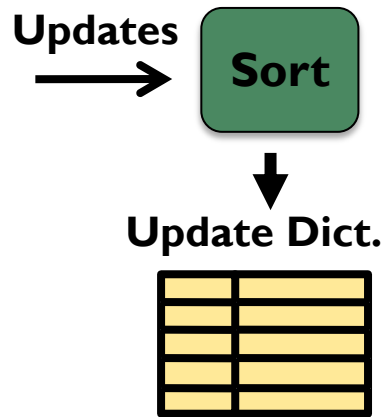
1 A simple tuple update in **row-wise layout** leads to **multiple random accesses** in **column-wise layout**

2 Updates change **encoded value** in the dictionary → (1) Need to **reconstruct** the dictionary, and (2) **recompress** the column

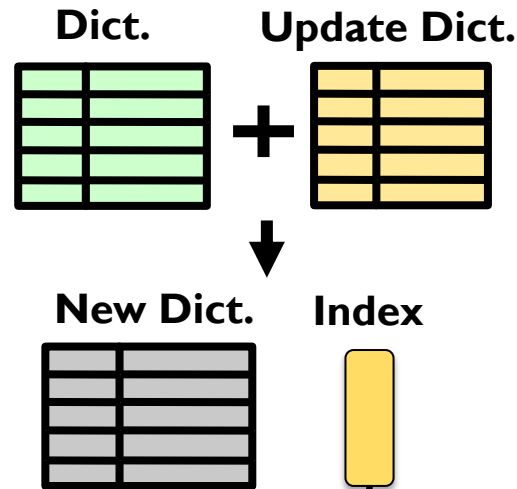
# Update Application: Algorithm

We design our update application algorithm to be aware of **PIM logic** characteristics and constraints

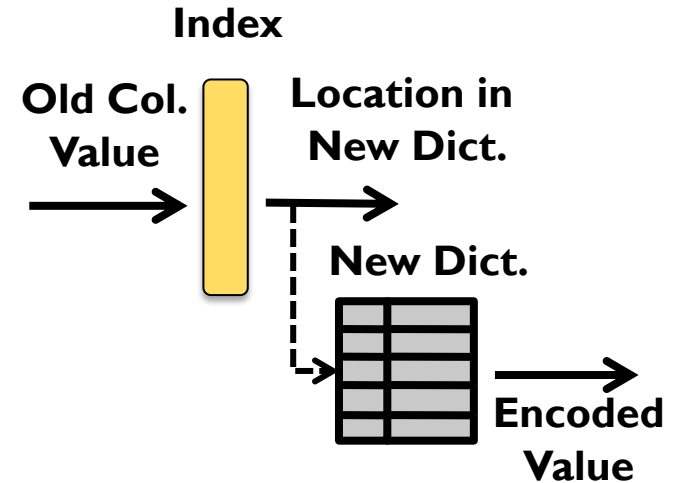
## Build Update Dict.



## Build New Dict. and Index



## New Compressed Col.

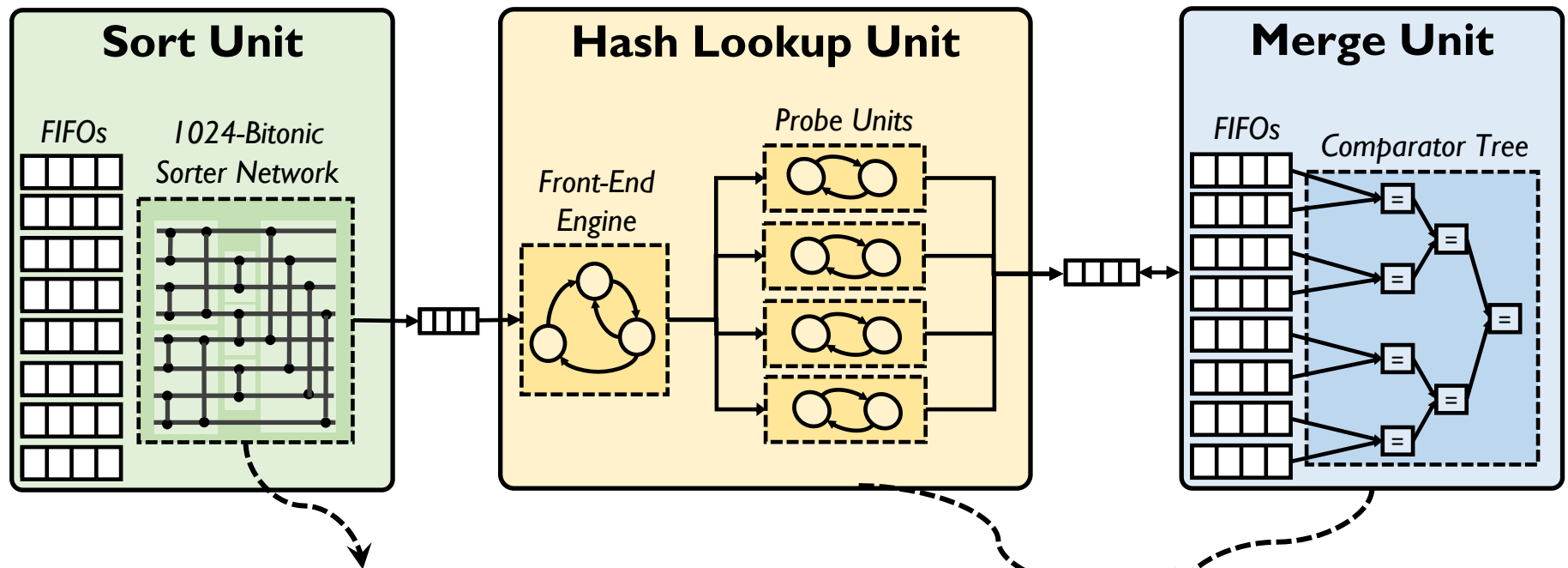


We maintain a **hash index** that links the **old encoded value** in a column to the **new encoded value**

Avoids the need to decompress the column and add updates, eliminating **data movement** and **random accesses** to 3D DRAM

# Update Application: Hardware

We design a **hardware implementation of our algorithm**, and add it to each **in-memory analytical island**



A **1024-value bitonic sorter**, whose basic building block is a network of comparators

Similar design as our **update gathering & shipping unit**

# Outline

1

Introduction

2

Limitations of HTAP Systems

3

Polynesia: Overview

4

Update Propagation Mechanism

5

**Consistency Mechanism**

6

Analytical Engine

7

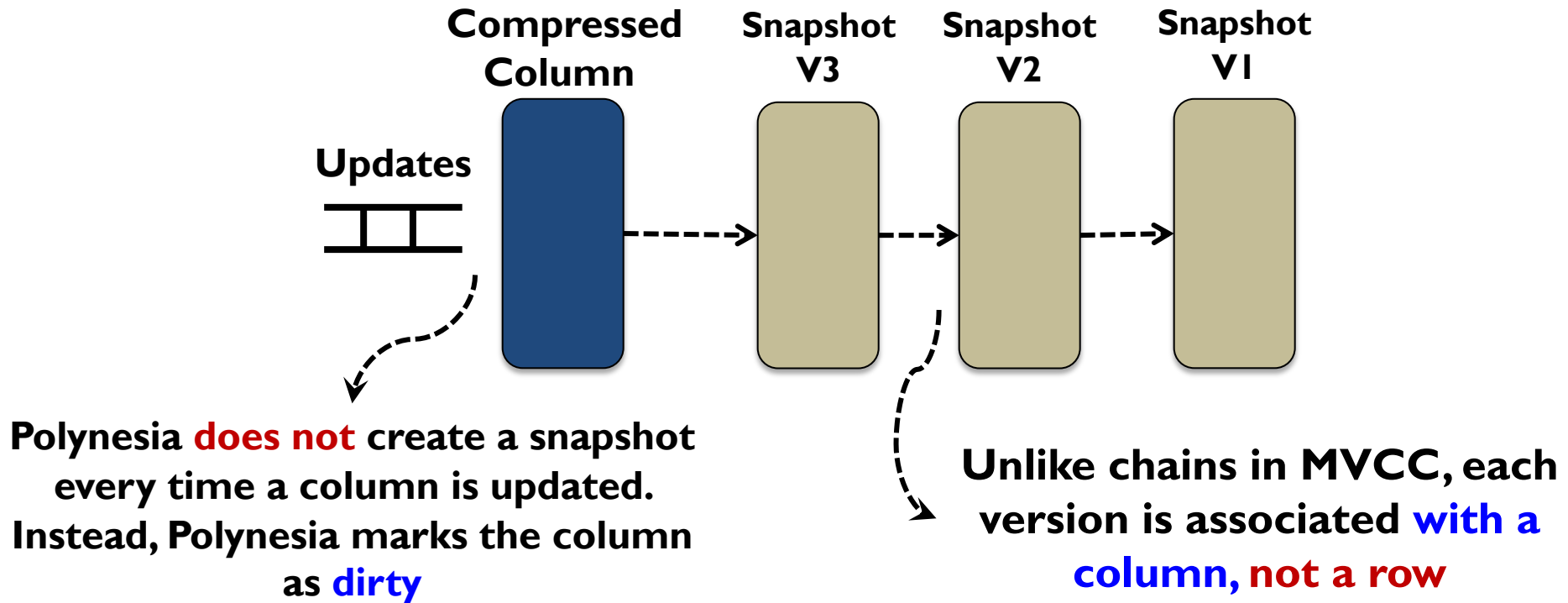
Evaluation

8

Conclusion

# Consistency Mechanism: Algorithm

For each column, there is **a chain of snapshots** where each chain entry corresponds to a **version of the column**

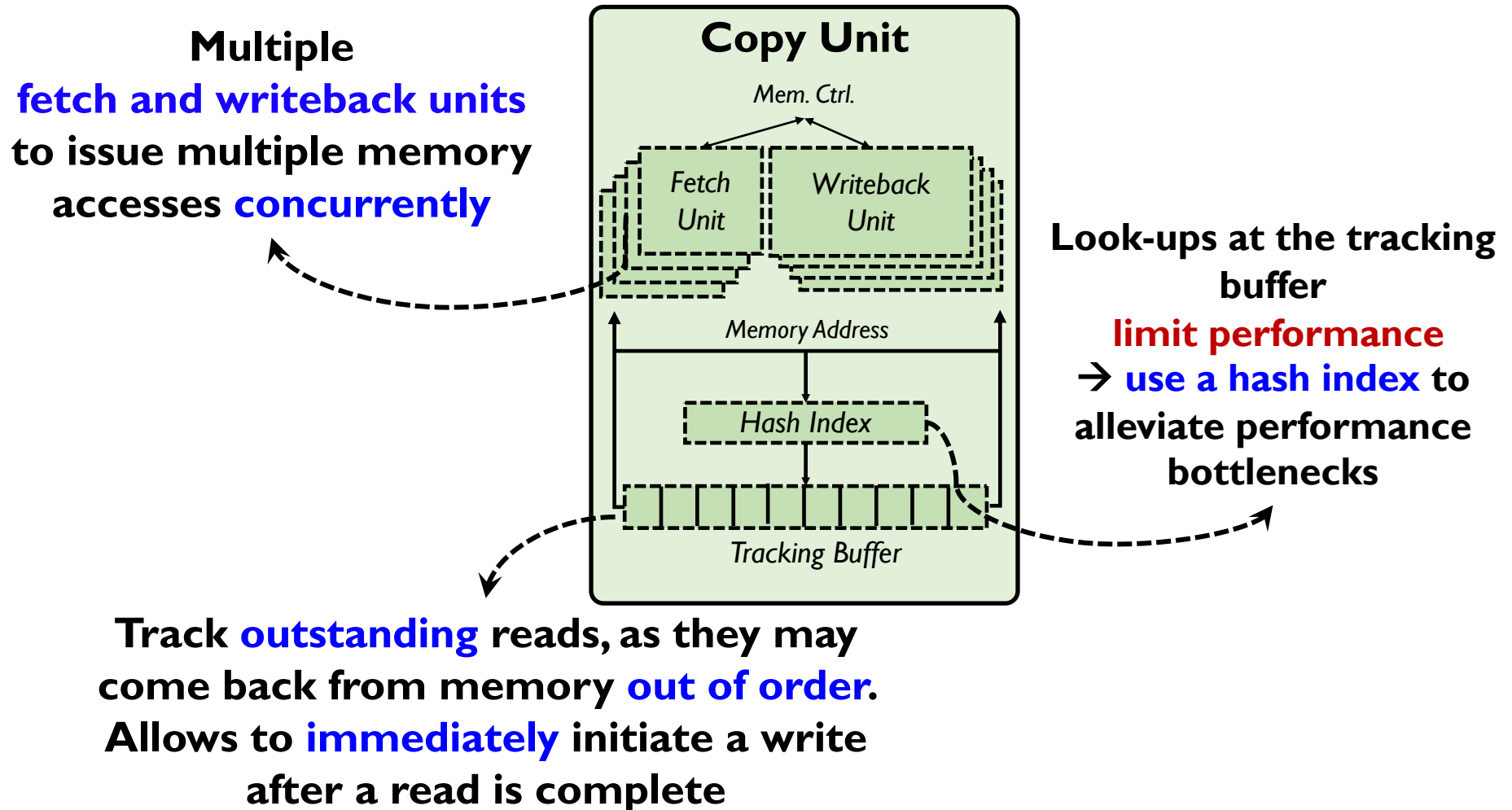


**Polynesia creates a new snapshot only if**

- (1) any of the columns are dirty, and**
- (2) no current snapshot exists for the same column**

# Consistency Mechanism: Hardware

Our algorithm success at satisfying **performance isolation** relies on how fast we can do **memcpy** to minimize **snapshotting latency**



# Outline

1

Introduction

2

Limitations of HTAP Systems

3

Polynesia: Overview

4

Update Propagation Mechanism

5

Consistency Mechanism

6

**Analytical Engine**

7

Evaluation

8

Conclusion

# Analytical Engine: Query Execution

Efficient analytical query execution **strongly depends** on:

1

Data layout and data placement

2

Task scheduling policy

3

How each physical operator is executed

The execution of **physical operators** of analytical queries significantly benefit from **PIM**

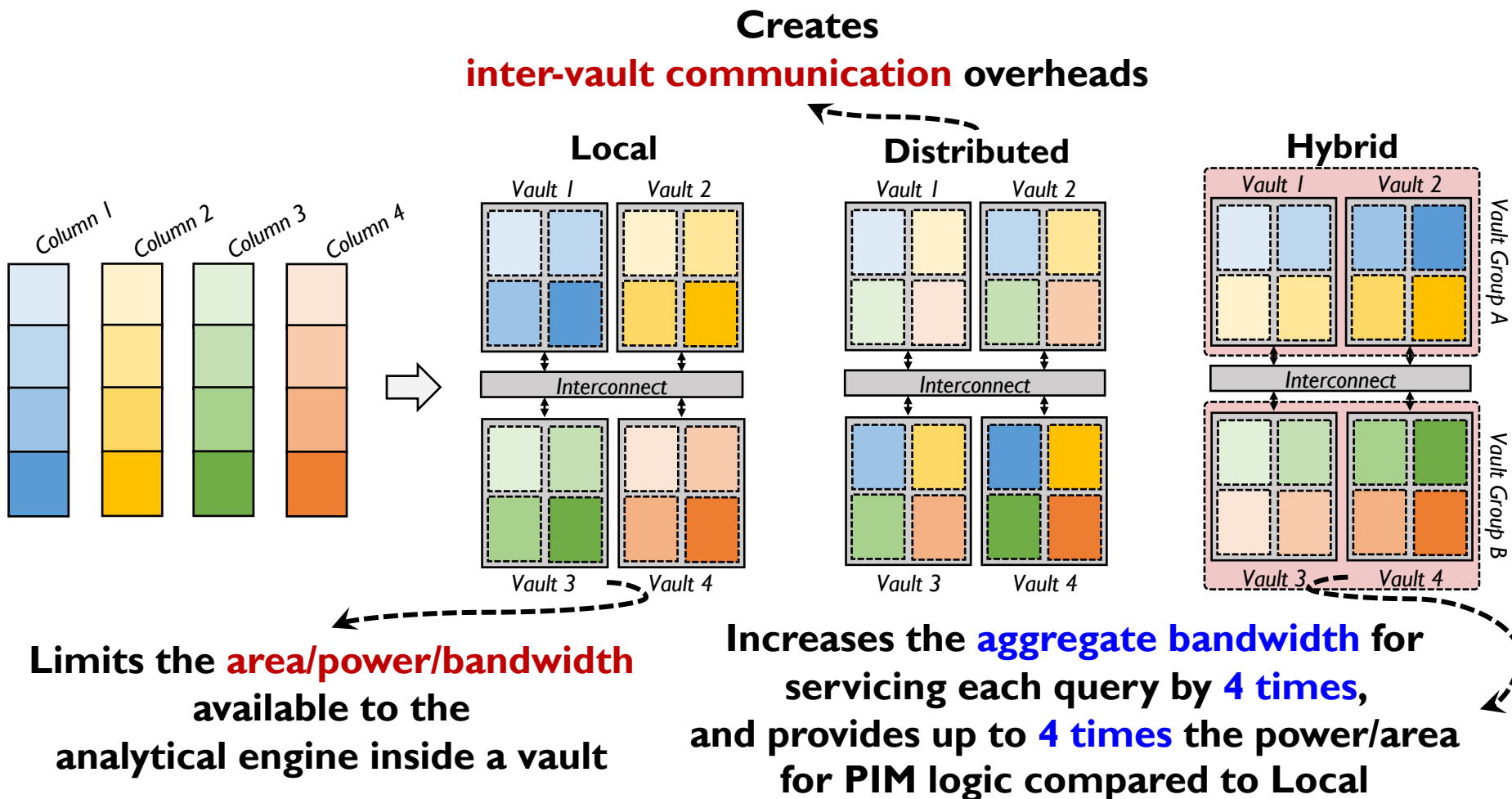


**Without PIM-aware data placement/task scheduler,  
PIM logic for operators alone cannot provide throughput**



# Analytical Engine: Data Placement

**Problem:** how to **partition analytical data** across vaults of the 3D-stacked memory



# Analytical Engine: Query Execution

Other details in the paper:

## Task scheduling policy

We design a **pull-based** task assignment strategy, where **PIM** threads **cooperatively** pull tasks from the task queue **at runtime**

## How each physical operator is executed

We employ the **top-down Volcano (Iterator)** execution model to execute physical operations (e.g., scan, filter, join) while respecting operator's dependencies

# Analytical Engine: Query Execution

Other details in the paper:

Task scheduling policy

## Polynesia: Enabling High-Performance and Energy-Efficient Hybrid Transactional/Analytical Databases with Hardware/Software Co-Design

Amirali Boroumand<sup>†</sup>  
<sup>†</sup>*Google*

Saugata Ghose<sup>◇</sup>  
<sup>◇</sup>*Univ. of Illinois Urbana-Champaign*

Geraldo F. Oliveira<sup>‡</sup>  
<sup>‡</sup>*ETH Zürich*

Onur Mutlu<sup>‡</sup>

We employ the **top-down Volcano (Iterator)** execution model to execute physical operations (e.g., scan, filter, join) while respecting operator's dependencies



Full Draft

# Outline

1

Introduction

2

Limitations of HTAP Systems

3

Polynesia: Overview

4

Update Propagation Mechanism

5

Consistency Mechanism

6

Analytical Engine

7

Evaluation

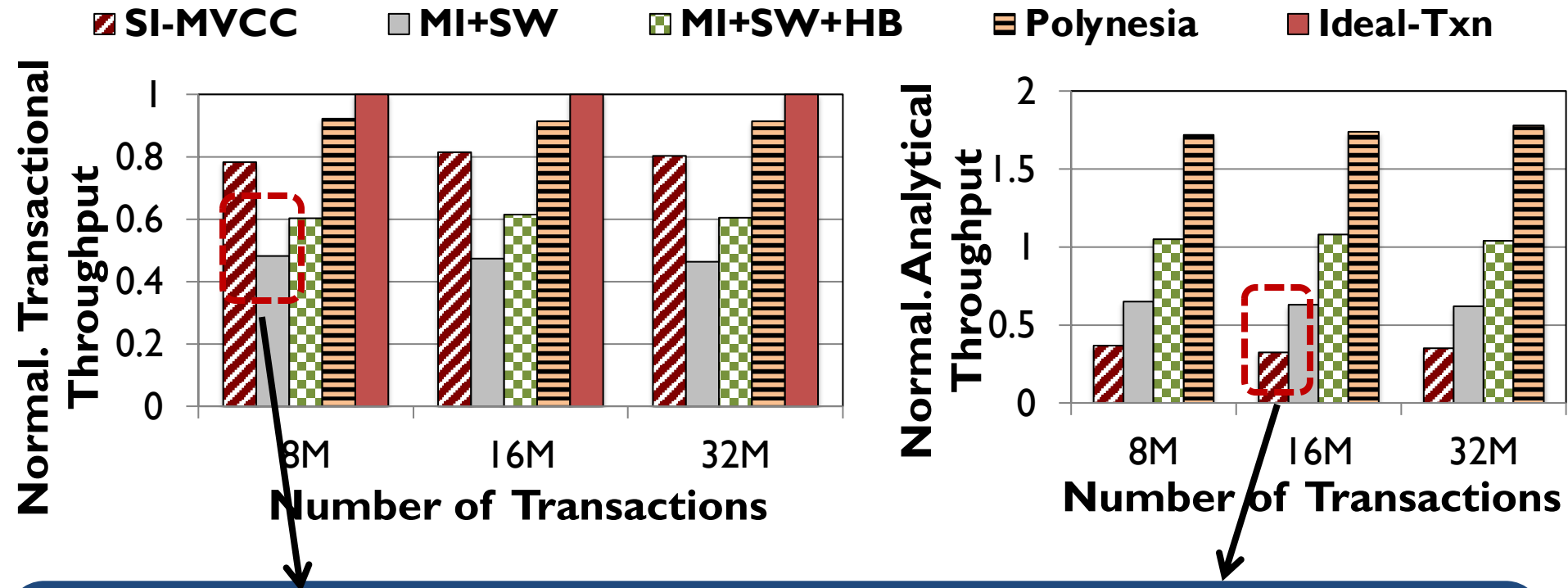
8

Conclusion

# Methodology

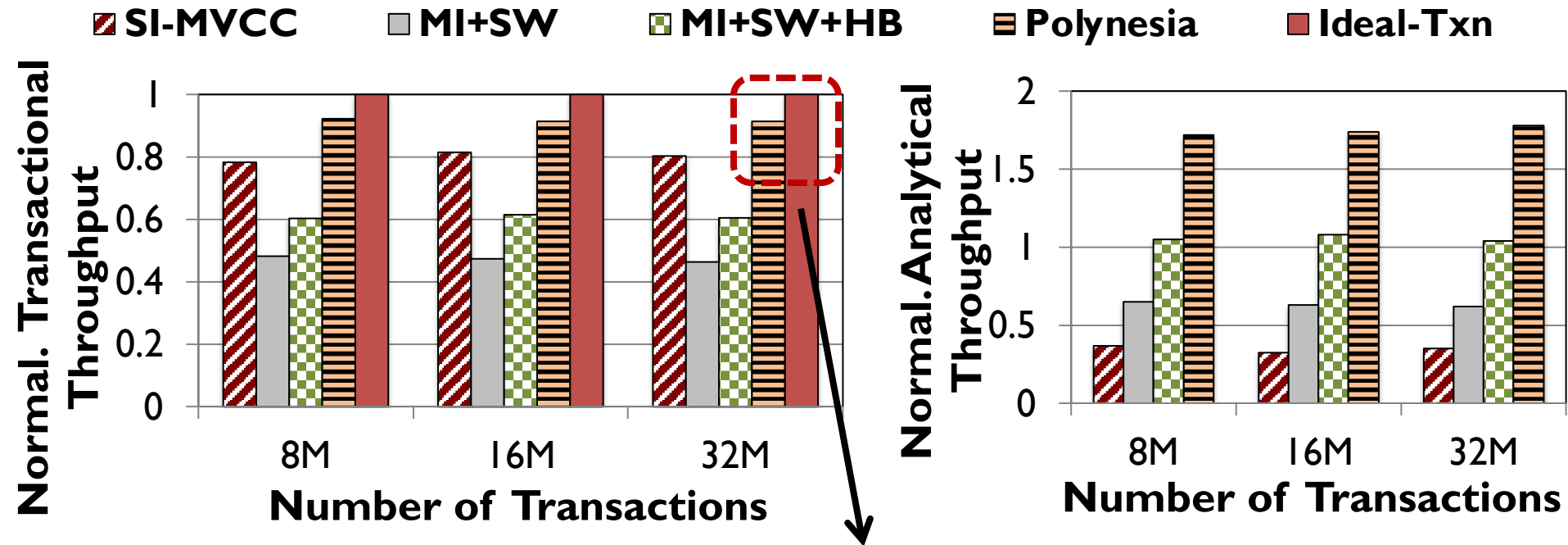
- We adapt previous transactional/analytical engines with our new algorithms
  - **DBx1000** for transactional engine
  - **C-store** for analytical engine
- We use **gem5** to simulate Polynesia
  - Available at: <https://github.com/CMU-SAFARI/Polynesia>
- We compare **Polynesia** against:
  - Single-Instance-Snapshotting (**SI-SI**)
  - Single-Instance-MVCC (**SI-MVCC**)
  - Multiple-Instance + Polynesia's new algorithms (**MI+SW**)
  - **MI+SW+HB**: MI+SW with a 256 GB/s main memory device
  - **Ideal-Txn**: the peak transactional throughput if transactional workloads run in isolation

# End-to-End System Analysis (1/5)

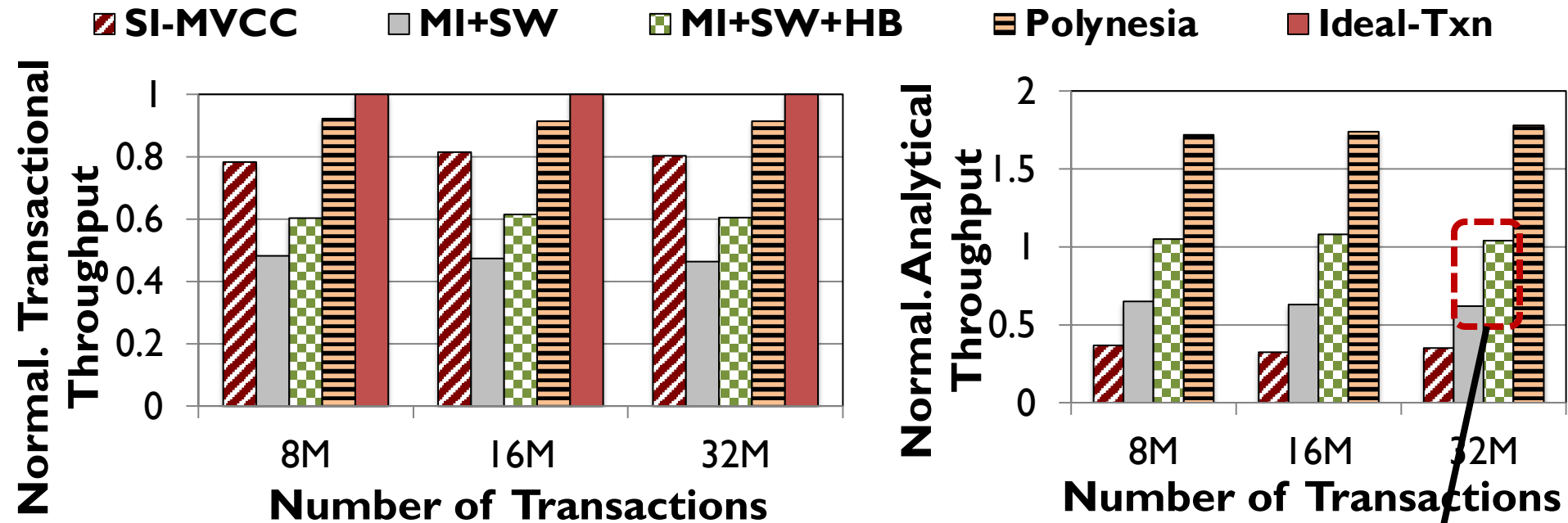


While SI-MVCC is the best baseline for **transactional throughput**, it degrades **analytical throughput** by **63.2%**, due to its **lack of workload-specific optimizations** and **consistency mechanism**

# End-to-End System Analysis (2/5)



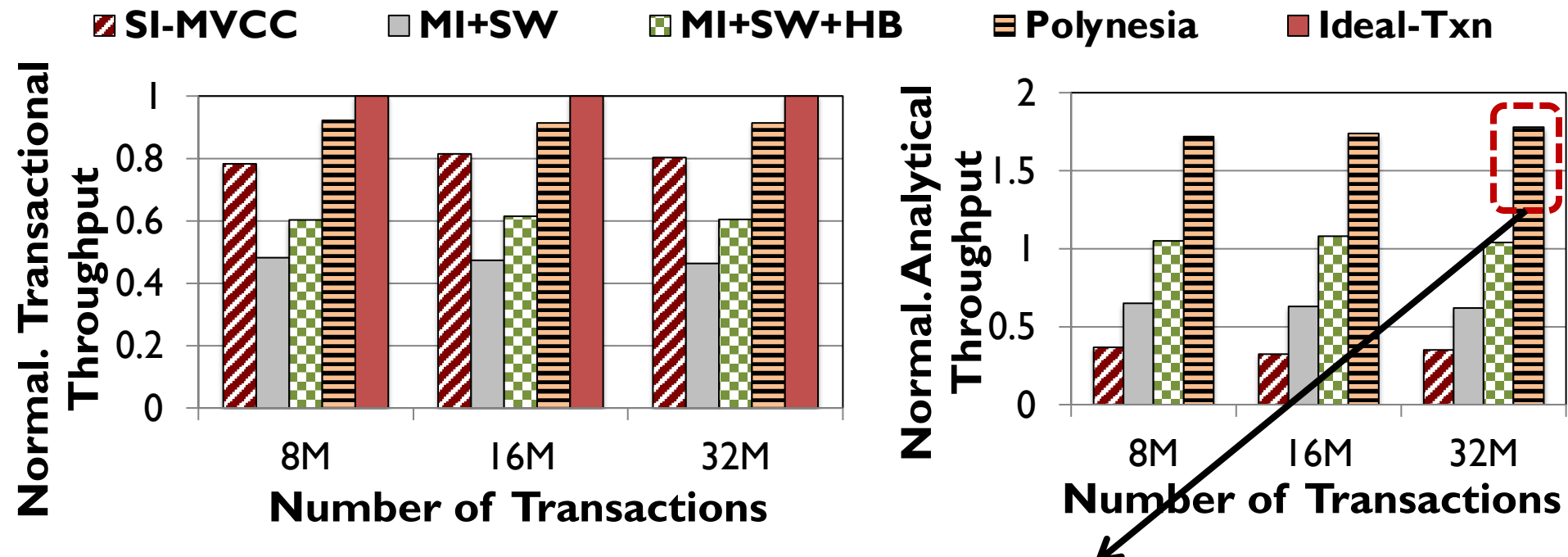
# End-to-End System Analysis (3/5)



**MI+SW+HB is the best software-only HTAP for analytical workloads, because it provides workload-specific optimizations, but it still loses 35.3% of the analytical throughput due to high main memory contention**

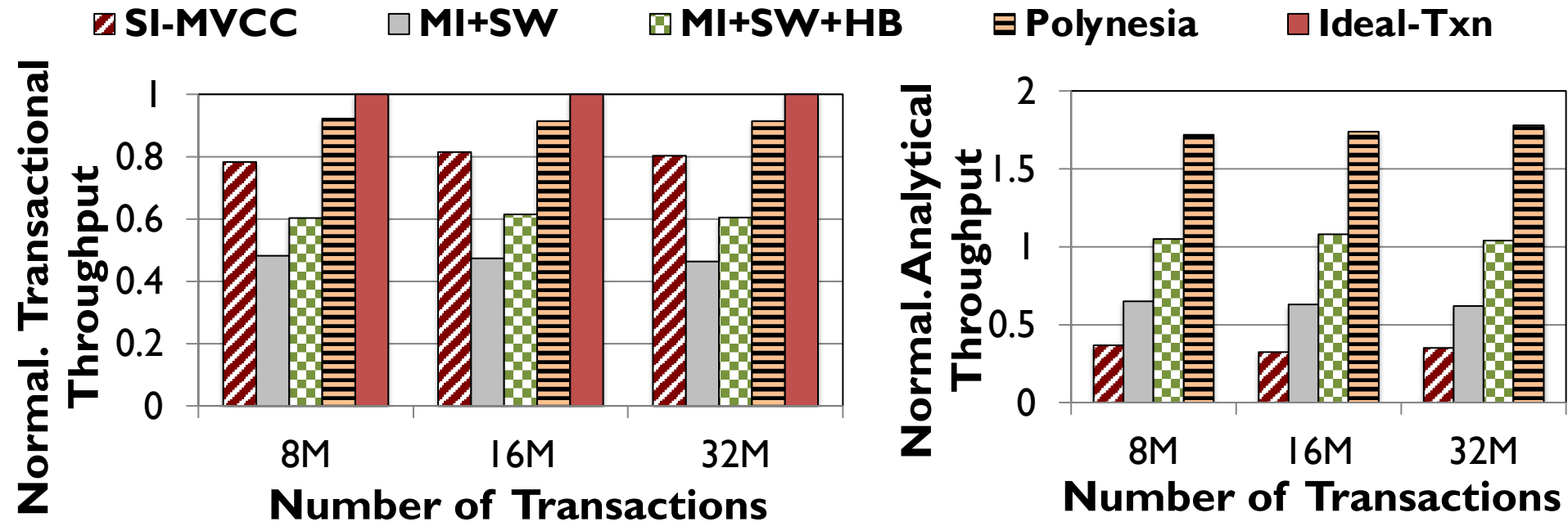


# End-to-End System Analysis (4/5)



Polynesia improves over **MI+SW+HB** by **63.8%**, by eliminating **data movement**, and using **custom logic** for **update propagation** and **consistency**

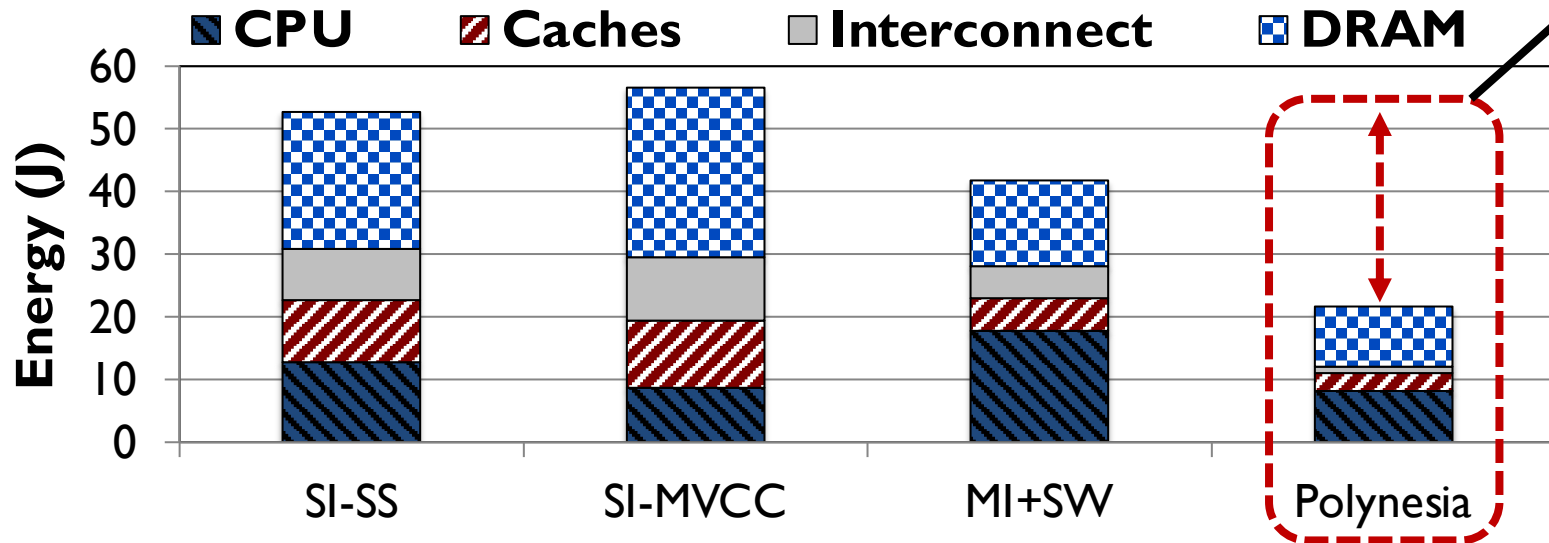
# End-to-End System Analysis (5/5)



Overall, Polynesia **achieves** all three **properties of HTAP** system and has a **higher** transactional/analytical **throughput (1.7x/3.74x)** over prior HTAP systems

# Energy Analysis

Polynesia consumes **0.4x/0.38x/0.5x** the energy of SI-SS/SI-MVCC/MI+SW since Polynesia **eliminates** a large fraction (**30%**) of **off-chip DRAM accesses**



Polynesia is an **energy-efficient HTAP system**,  
**reducing** energy consumption by **48%**,  
on average across prior works

# More in the Paper

- Real workload analysis
- Effect of the update propagation technique
- Effect of the consistency mechanism
- Effect of the analytical engine
- Effect of the dataset size
- Area Analysis

# More in the Paper

- Real workload analysis

- Effect of the update propagation technique

## Polynesia: Enabling High-Performance and Energy-Efficient Hybrid Transactional/Analytical Databases with Hardware/Software Co-Design

Amirali Boroumand<sup>†</sup>  
<sup>†</sup>*Google*

Saugata Ghose<sup>◇</sup>  
<sup>◇</sup>*Univ. of Illinois Urbana-Champaign*

Geraldo F. Oliveira<sup>‡</sup>  
<sup>‡</sup>*ETH Zürich*

Onur Mutlu<sup>‡</sup>

- Effect of the dataset size

- Area Analysis



Full Draft

# Outline

1

Introduction

2

Limitations of HTAP Systems

3

Polynesia: Overview

4

Update Propagation Mechanism

5

Consistency Mechanism

6

Analytical Engine

7

Evaluation

8

Conclusion

# Conclusion

- **Context:** Many applications need to perform real-time data analysis using an Hybrid Transactional/Analytical Processing (HTAP) system
  - An ideal HTAP system should have **three properties**:  
(1) **data freshness** and **consistency**, (2) **workload-specific optimization**,  
(3) **performance isolation**
- **Problem:** Prior works **cannot achieve all properties** of an ideal HTAP system
- **Key Idea:** Divide the system into transactional and analytical **processing islands**
  - Enables **workload-specific optimizations** and **performance isolation**
- **Key Mechanism:** Polynesia, a novel hardware/software cooperative design for in-memory HTAP databases
  - Implements **custom algorithms and hardware** to reduce the costs of **data freshness** and **consistency**
  - Exploits **PIM** for analytical processing to alleviate **data movement**
- **Key Results:** Polynesia outperforms three state-of-the-art HTAP systems
  - Average transactional/analytical throughput improvements of **1.7x/3.7x**
  - **48%** reduction on energy consumption

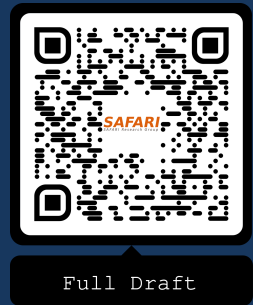
# Polynesia:

## Enabling High-Performance and Energy-Efficient Hybrid Transactional/Analytical Databases with Hardware/Software Co-Design

P&S Processing-in-Memory

Fall 2022

17 January 2023



**Amirali Boroumand**  
**Geraldo F. Oliveira**

**Saugata Ghose**  
**Onur Mutlu**



# Heterogeneous Data-Centric Architectures for Modern Data-Intensive Applications: Case Studies in Machine Learning and Databases

P&S Processing-in-Memory  
Fall 2022  
17 January 2023

**Geraldo F. Oliveira**

**SAFARI**

**ETH** zürich